

Bedienung einer Getränkemischmaschine über Sprachbefehle

STUDIENARBEIT

für die Prüfung zum

Bachelor of Science

des Studienganges Angewandte Informatik

an der

Dualen Hochschule Baden-Württemberg Karlsruhe

von

Felix Manuel Gervasi

und

Alena Sutiagina

Abgabedatum 31. März 2023

Bearbeitungszeitraum	6 Monate
Matrikelnummer	1052491
Kurs	TINF20B4
Ausbildungsfirma	

Betreuer der Studienarbeit	Prof. Dr. Jörn Eisenbiegler
----------------------------	-----------------------------

Erklärung

Ich versichere hiermit, dass ich meine Studienarbeit mit dem Thema: »Bedienung einer Getränkemischmaschine über Sprachbefehle« selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Ich versichere zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

Ort Datum

Unterschrift

Zusammenfassung

TODO

Abstract

TODO

Inhaltsverzeichnis

Abbildungsverzeichnis	iii
Tabellenverzeichnis	iii
Listingverzeichnis	v
Formelverzeichnis	v
Akürzungsverzeichnis	vii
1 Einleitung	1
1.1 Aufgabenstellung	1
1.2 Vorgehen	2
2 Theoretische und technische Grundlagen	3
2.1 Getränkemischmaschine	3
2.2 Hardware	4
2.3 Sprachverarbeitung	4
2.3.1 Verarbeitung natürlicher Sprache	6
2.3.2 Tokenisierung von Wörtern	8
2.3.3 Vektorisierung von Wörtern	10
2.3.4 Syntaktische Analyse	21
2.3.5 Ansätze für die Erstellung eines Chatbots	27
2.3.6 Mehrschichtiges Perzeptron für ein Mehrklassen-Klassifizierungsproblem	33
3 Anforderungen	36
3.1 Antwortzeit	36
3.2 Offline-Funktionalität	36
3.3 Lautstärke	37
3.4 Entfernung	37
3.5 Antworten	37

3.6	Kosten	37
3.7	Verbrauch von Arbeits- und Festplattenspeicher	38
3.8	Anpassungsfähigkeit	38
3.9	Bewertung der Ansätze zur Erstellung eines Dialogsystems	39
4	Konzept	40
4.1	Allgemein	40
4.2	Bewertungskriterien	40
4.3	Konzept A: Spracherkennung und -verarbeitung mittels Arduino	41
4.4	Konzept B: Spracherkennung und -verarbeitung mittels mobiler Anwendung	43
4.5	Konzept C: Spracherkennung und -verarbeitung auf Computer-Hardware	44
4.6	Finales Hardware-Konzept	45
4.7	Konzept für die Sprachsteuerung	46
4.7.1	Ansatz für das Dialogsystem	46
4.7.2	Befehle	48
5	Implementierung	49
5.1	Implementierung des Sprachverarbeitungssystems	49
5.1.1	Word2Vec-Modell	49
5.1.2	Sequence-to-Sequence-Modell	49
5.2	Implementierung der Sprachsteuerung	49
5.2.1	Spracherkennung	49
5.2.2	Anbindung des Sprachmodells an die Mischmaschine	55
5.2.3	Befehlsverarbeitung in der Mischmaschine	55
6	Fazit und Ausblick	59
	Literaturverzeichnis	ix
	Anhang	ix

Abbildungsverzeichnis

2.1	Schematischer Aufbau der Bedienungsschnittstelle	4
2.2	Rekurrente Chat-Bot-Pipeline	7
2.3	Unitärer Vektor	12
2.4	Normalisierte Term Frequency (TF)-Vektoren	13
2.5	Normalisierte TF-Vektoren mit Python	14
2.6	Zweidimensionale Vektoren von Termhäufigkeiten	15
2.7	Euklidische Norm	15
2.8	Zweidimensionale Theta	16
2.9	Euklidisches Skalarprodukt	16
2.10	Ochiai-Koeffizient	16
2.11	Singulärwertzerlegung (SWZ)-Matrize	24
2.12	Trunkierte SWZ-Matrize	25
2.13	Artificial Intelligence Markup Language (AIML) Chatbot	28
2.14	Mehrschichtiges Perzeptron	34
4.1	Spracherkennung und -verarbeitung mittels Arduino	42
4.2	Spracherkennung und -verarbeitung mittels mobiler Anwendung	43
4.3	Spracherkennung und -verarbeitung auf Computer-Hardware	44
5.1	Fehlermeldungen beim Zugriff auf das Mikrofon	54

Tabellenverzeichnis

2.1	TF-Vektoren	19
2.2	Inverse Document Frequency (IDF)-Vektoren	19
2.3	Term Frequency-Inverse Document Frequency (TF-IDF)-Vektoren	20
3.1	Bewertung der Ansätze zur Erstellung eines Dialodsystems	39
4.1	Bewertung der Konzepte	45
4.2	Bewertung der Ansätze für die Erstellung eines Dialogsystems	47

Listingverzeichnis

5.1	Audioaufnahme mit <i>SpeechRecognition</i>	50
5.2	Sprache zu Text mit <i>OpenAI Whisper</i>	52
5.3	Hot-Word-Erkennung	52
5.4	Fehlerbehandlung	55
5.5	Serielle Kommunikation in Python	55
5.6	Serielle Kommunikation in Python	56
5.7	Serielle Kommunikation für die Mischmaschine	57
5.8	Interpretation der Kommandos	58

Formelverzeichnis

(2.1) Wert der Worthäufigkeit im Dokument	18
(2.2) Wert der inversen Dokumenthäufigkeit	18
(2.3) TF-IDF-Wert	18
(2.4) Singularitätwertzerlegung	23
(2.5) Trunkierte SWZ-Matrix	25
(2.6) Entropie H	35

Abkürzungsverzeichnis

HTTP Hypertext Transfer Protocol	41
NLP Natural Language Processing	5
AIML Artificial Intelligence Markup Language	iii
XML eXtensible Markup Language	28
NLTK Natural Language Toolkit	8
TF Term Frequency	iii
IDF Inverse Document Frequency	iv
RBM Restricted Boltzmann Machines	32
GAN Generative Adversarial Network	32
RNN Recurrent Neural Networks	33
IoT Internet der Dinge	1
KI Künstlichen Intelligenz	2
ML Machine Learning	2
UI User Interface	4
RAM Read Only Memory	38
CFG Context-free grammar	8
NLTK Natural Language Toolkit	8
API Application Programming Interface	8
CMU Carnegie Mellon University	50
BOW Bag Of Words	11
seq2seq sequence-to-sequence	12
GB Gigabyte	54
TF-IDF Term Frequency-Inverse Document Frequency	iv
MLP Multilayer Perceptron	33

LSA Latent Semantic Analysis	21
SWZ Singulärwertzerlegung	iii
SVD Singular Value Decomposition	22
MLP Mehrschichtiges Perzeptron	33

Kapitel 1

Einleitung

Die Informationstechnik versteckt sich heutzutage fast überall - selbst dort, wo sie von den meisten Menschen nicht vermutet werden würde. Beispiele hierfür sind Autos, Kaffeemaschinen, Zahnbürsten, Rasierer, Küchengeräte und vieles mehr. Grund dafür ist die fortschreitende Möglichkeit der Miniaturisierung von Computern, sodass diese nahezu überall verbaut werden können. Beispielsweise können Mikrochips in der Kaffeemaschine dafür sorgen, dass die richtige Menge an Kaffee serviert wird oder der Füllstand der einzelnen Behälter angezeigt werden kann. Solche Systeme, die Informationen mit Hilfe eines Computers verarbeiten und dabei mit ihrer Umgebung derartig „verschmelzen“, nennt man auch *embedded systems* (z. Dt. *eingebettete Systeme*) [MARWEDEL 2021].

Der aktuelle Trend des Internet der Dinge (IoT) führt zu einem noch größeren Anstieg eingebetteter Systeme im Alltag. Im IoT geht es speziell um eingebettete Systeme, die internetfähig (vernetzt) sind. Nach Schätzungen des Marktforschungsunternehmens *Gartner* gab es im Jahr 2017 8,4 Milliarden solcher vernetzten Geräte weltweit [JANSEN 2017]. Das die Menge der vernetzten Geräte als Teilmenge der eingebetteten Systeme betrachtet werden kann ist damit zu rechnen, dass deren Anzahl sogar weit größer ausfällt.

1.1 Aufgabenstellung

Im Rahmen dieser Arbeit geht es um die Sprachsteuerung einer Getränkemischmaschine, die in diesem Fall als eingebettetes System zu verstehen ist und in einer vorangegangenen Arbeit bereits konzipiert und gebaut wurde [GÖTZ, HÖCKELE und LOBERT 2022]. Sie verfügt derzeit über ein Touch-Display zur Bedienung durch den Benutzer. Ziel der Arbeit ist es zusätzlich eine natürlischsprachliche Interaktion mit der Maschine zu ermöglichen, die mindestens den Funktionsumfang besitzt, der aktuell über die Bildschirmeingabe möglich ist. Dabei soll die Maschine nicht nur in der Lage sein die natrliche Sprache des Benutzers in ein geeignetes Format umzuwandeln, sodass die Maschine den korrekten Befehl ausführt. Sie soll auch in der

Lage sein dem Benutzer zu Antworten, sodass die Illusion einer Konversation mit der Maschine entsteht.

1.2 Vorgehen

Zunächst müssen die Sprachverarbeitung und Spracherkennung betrachtet werden. Die Sprachverarbeitung dient der Interpretation des Gesprochenen um eine geeignete Antwort auszugeben und dem Übersetzen in einen Maschinenbefehl. Im Rahmen dieser Arbeit sollen dafür Verfahren und Techniken der Künstlichen Intelligenz (KI) und des Machine Learning (ML) eingesetzt werden. Die Spracherkennung beschäftigt sich mit der Aufnahme des Tonsignals bzw. der Schallwellen (bspw. über ein Mikrofon) und dem Umwandeln dieser Signale in Text, sodass dieser an das KI-Modell weitergereicht werden kann.

Bei der Arbeit mit eingebetteten Systemen muss man sich der vorhandenen Hardwareleistung und den benötigten Hardwareanforderungen bewusst sein, da diese meist sehr begrenzt ist. Deshalb werden im Rahmen dieser Arbeit verschiedene Ansätze diskutiert, wie und wo die einzelnen Schritte und Berechnungen ablaufen sollen (s. Kapitel 4).

Kapitel 2

Theoretische und technische Grundlagen

2.1 Getränkemischmaschine

Ziel dieser Arbeit ist die Implementierung einer Sprachsteuerung für eine Getränkemischmaschine. Die Getränkemischmaschine wurde bereits in einem vorangegangenen Projekt erstellt [GÖTZ, HÖCKELE und LOBERT 2022]. In diesem Abschnitt wird darauf eingegangen, um was für eine Art von Maschine es sich dabei handelt und es werden ihre Funktionsweise und ihr Aufbau beschrieben.

Das Mischen von Getränken bzw. Flüssigkeiten ist der Anwendungsfall für den die Maschine konzipiert wurde. Dazu besitzt die Mischmaschine fünf Behälter zu je einem Liter. Jedem Behälter ist eine Pumpe zugeordnet, die separat angesteuert werden kann und dafür sorgt, dass die Flüssigkeit aus dem Behälter zur Getränkeausgabe gelangt. Kurz vor dem Ausgang werden die Schläuche der Behälter zusammengeführt, wodurch letztlich die Mischung der verschiedenen Getränke erzielt wird. Sogenannte „Rückschlagventile“ sorgen dafür, dass ein Zurückfließen der Getränkemischung in die Behälter verhindert wird.

Zur Steuerung der Maschine durch einen Benutzer befinden sich an der Vorderseite zwei Knöpfe und ein Touch-Display. Einer der Knöpfe dient dem Anschalten der Maschine und ein weiterer der Ausgabe des Getränks. Über das Touch-Display kann der Benutzer die Mischung seiner Getränke konfigurieren und die Maschine administrieren.

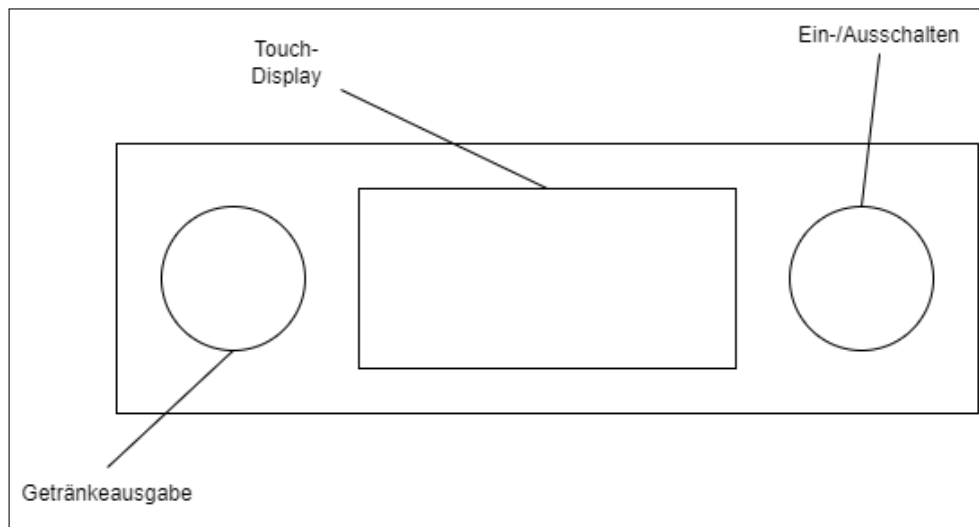


Abbildung 2.1: Schematischer Aufbau der Bedienungsschnittstelle

Abbildung 2.1 stellt den schematischen Aufbau der Bedienungsschnittstelle der Getränkemischmaschine für ein besseres Verständnis dar.

Die technische Umsetzung basiert auf der Kommunikation zwischen dem Touch-Display und einem, in der Getränkemischmaschine verbauten, Arduino, der anhand der Daten von der Benutzereingabe, die Pumpen steuert. Beim Drücken des Startknopfes werden das Display und der Arduino mit Strom versorgt, sodass sowohl das Display als auch der Arduino starten und mit der Ausführung des benutzerdefinierten Quelltextes beginnen. Auf dem Bediendisplay findet sich fünf Schieberegler - ein Schieberegler je Behälter - mit denen der Benutzer die Zusammensetzung seines Mischgetränks aus den fünf Behältern konfigurieren kann. Beim Drücken eines in der User Interface (UI) des Displays dargestellten Knopfes werden die Werte der Schieberegler an den Arduino übertragen. Dieser rechnet die Prozentwerte der Schieberegler in Durchsatzraten für die Pumpen um. Wird anschließend der Knopf für die Getränkeausgabe gedrückt gehalten steuert der Arduino die Pumpen mit ihrer jeweiligen Durchsatzrate an und der Anwender bekommt sein Mischgetränk ausgegeben.

2.2 Hardware

2.3 Sprachverarbeitung

Im Zuge des technologischen Fortschritts nutzen die Menschen heutzutage zunehmend Sprachassistenten für verschiedene Aufgaben. Einer der Hauptvorteile der Sprachsteuerung ist die Bequemlichkeit und Geschwindigkeit, mit der Aufgaben erledigt werden können, ohne dass man tippen oder mit der Maus klicken muss. Sprachassistenten nutzen die Verarbeitung natürlicher

Sprache, um Befehle zu erkennen und zu verstehen, die der Nutzer laut ausspricht.

Die Verarbeitung natürlicher Sprache - Natural Language Processing (NLP) - ist eine wichtige Technologie, die es Computern ermöglicht, die von Menschen verwendete natürliche Sprache zu verstehen. Diese Technologie ermöglicht es Computern, menschliche Sprache zu erkennen und zu interpretieren und Text und Sprache in natürlicher Sprache zu erzeugen. Im Bereich der Sprachsteuerung spielt NLP eine Schlüsselrolle bei der Erkennung von Sprache und dem Verstehen von Befehlen, die der Benutzer laut ausspricht. Es wird verwendet, um die Sprache des Benutzers in Text umzuwandeln, den ein Computer verstehen und verarbeiten kann. Um dies zu erreichen, verwendet NLP eine Vielzahl von Techniken und Technologien, darunter maschinelles Lernen, Tonanalyse, syntaktische Analyse und mehr. [JURAFSKY u. a. 2009]

Im Rahmen dieses Projekts erfordert die Implementierung der Sprachsteuerung einer Getränkemischmaschine die Verarbeitung natürlicher Sprache, damit die Maschine Befehle verstehen kann, die der Benutzer laut ausspricht. Dieses Projekt ähnelt einem Chatbot, bei dem der Benutzer eine Frage stellen oder einen Befehl geben kann und der Chatbot führt die entsprechende Aktion aus. Die Verarbeitung natürlicher Sprache ist für die Entwicklung eines solchen Sprachsteuerungssystems unerlässlich und ermöglicht es der Maschine, Befehle in natürlicher Sprache zu verstehen und auszuführen.

Eine der wichtigsten Komponenten der Verarbeitung natürlicher Sprache ist die Spracherkennung und das Syntaxanalyseverfahren. Bei der Spracherkennung kommen Deep-Learning-Techniken zum Einsatz, die es einem Computer ermöglichen, Sprachlaute zu erkennen und in Text zu übersetzen. Anschließend wird das Syntaxanalyseverfahren verwendet, um die Satzstruktur zu bestimmen und Schlüsselwörter und -sätze hervorzuheben, die zur Bestimmung des Benutzerbefehls verwendet werden können. Auch die Tonwertanalyse ist ein wichtiger Bestandteil der Verarbeitung natürlicher Sprache. Mit Hilfe der Tonalitätsanalyse lässt sich die emotionale Färbung des Textes bestimmen, was für die Ermittlung der Absicht des Nutzers nützlich sein kann. Da die Tonwertanalyse im Rahmen dieser Arbeit nicht relevant ist, wird sie nicht weiter erörtert.

Darüber hinaus werden grammatik- und regelbasierte Technologien zur Verarbeitung natürlicher Sprache eingesetzt. Diese Technologien werden eingesetzt, um die korrekte Struktur des Benutzerbefehls zu bestimmen und Schlüsselwörter hervorzuheben, die zur Durchführung von Aktionen verwendet werden können. Außerdem werden Techniken des maschinellen Lernens eingesetzt, damit der Computer aus früheren Befehlen und Aktionen „lernen“ kann, was die Genauigkeit der Erkennung von Benutzerbefehlen verbessert.

Daher ist der Einsatz von Technologien zur Verarbeitung natürlicher Sprache für das Projekt der Sprachmischmaschine unerlässlich. Die Verarbeitung natürlicher Sprache wird es der Maschine ermöglichen, die Befehle zu verstehen und zu verarbeiten, die der Benutzer laut ausspricht und die entsprechenden Aktionen durchzuführen.

2.3.1 Verarbeitung natürlicher Sprache

Die Verarbeitung natürlicher Sprache ist ein Forschungsbereich der Informatik und der KI, der sich mit der Verarbeitung natürlicher Sprache befasst. Bei der Verarbeitung geht es in der Regel darum, natürliche Sprache in Daten zu übersetzen, die ein Computer nutzen kann, um Informationen über die Welt um ihn herum zu erhalten.

Die NLP-Pipeline, die zur Erstellung eines Dialogsystems erforderlich ist, erfordert vier Arten der Verarbeitung sowie eine Datenbank zur Speicherung vergangener Äußerungen und Antworten. Jeder dieser Schritte kann einen oder mehrere Verarbeitungsalgorithmen enthalten, die parallel oder sequentiell arbeiten [LANE, HOWARD und HAPKE 2019]:

- Syntaktische Zergliederung - Extraktion von Merkmalen (strukturierte numerische Daten) aus natürlichem Text.
- Analyse - Erstellen und Kombinieren von Items, um Indikatoren für den Ton, die grammatikalische Korrektheit und die Semantik des Textes zu erhalten.
- Generierung - Generierung möglicher Antworten mit Hilfe von Mustern, Suchwerkzeugen oder Sprachmodellen.
- Ausführung - bereitet Aussagen vor, die auf der Geschichte und dem Zweck des Gesprächs basieren und wählt eine Folgeantwort.

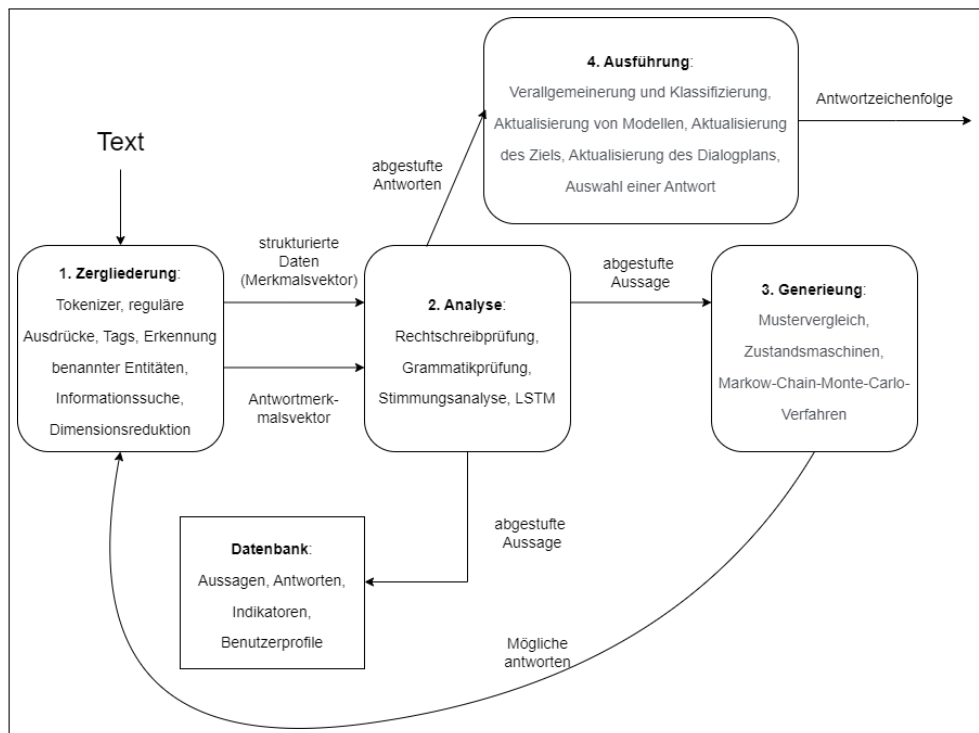


Abbildung 2.2: Rekurrente Chat-Bot-Pipeline

Die meisten Chatbots enthalten Elemente aus allen fünf Teilsystemen (die vier Verarbeitungsstufen sowie die Datenbank). Viele Anwendungen erfordern jedoch nur einfache Algorithmen, um viele dieser Schritte auszuführen. Ein Chatbot oder virtueller Assistent für Verbraucher wie Alexa oder Allo ist in der Regel so konzipiert, dass er äußerst sachkundig und leistungsfähig ist. Die Logik, die zur Beantwortung von Anfragen verwendet wird, ist jedoch oft oberflächlich und besteht aus einer Reihe von Codephrasen, die mit einer einzigen if-then-Entscheidungsverzweigung zur gleichen Antwort führen. Alexa (und die zugrundeliegende Lex-Engine) verhält sich wie ein einschichtiger, flacher Baum von Operatoren (if, elif, elif...). Andererseits stützt sich die Google Translate-Pipeline (oder jedes ähnliche maschinelle Übersetzungssystem) auf eine mehrstufige Hierarchie von Merkmalsextraktoren, Entscheidungsbäumen und Wissensgraphen, die Fragmente von Wissen über die Welt miteinander verbinden.

Alle NLP-Merkmale müssen gut funktionieren, damit das Dialogsystem richtig funktioniert:

- Merkmalsextraktion (normalerweise zur Erstellung eines Vektorraummodells).
- Informationssuche zur Beantwortung von Sachfragen.
- Semantische Suche, um Informationen aus zuvor aufgezeichneten natürlichsprachlichen Texten oder Dialogen zu assimilieren.
- Generierung von natürlicher Sprache, um neue sinnvolle Aussagen zu verfassen.

2.3.2 Tokenisierung von Wörtern

Eine der wichtigsten Aufgaben der Verarbeitung natürlicher Sprache ist die Tokenisierung, d. h. die Zerlegung von Text in einzelne Wörter oder Phrasen, die so genannten Token. Dieser Prozess ist ein wichtiger Bestandteil der Verarbeitung natürlicher Sprache, da er es Programmen und Algorithmen ermöglicht, Texte anhand ihrer Bestandteile zu verstehen und zu analysieren. Die Tokenisierung ist für die Verarbeitung natürlicher Sprache von entscheidender Bedeutung, da sie es Programmen ermöglicht, Schlüsselwörter, Phrasen und semantische Einheiten in einem Text hervorzuheben. Die Tokenisierung kann auch dazu verwendet werden, nicht benötigte Textelemente wie Satzzeichen, Leerzeichen und andere Zeichen, die keine Bedeutung haben, zu entfernen. Dadurch können Programme und Algorithmen effizienter und genauer arbeiten und die Wahrscheinlichkeit von Fehlern und falscher Textverarbeitung verringern.

Im NLP ist die Tokenisierung eine besondere Art der Segmentierung von Dokumenten. Bei der Segmentierung wird der Text in kleinere Abschnitte (Segmente) mit engerem Informationsgehalt unterteilt. Die Segmentierung kann die Aufteilung eines Dokuments in Absätze, in Sätze, in Phrasen und in Token (Wörter) sowie Satzzeichen umfassen. Ein Tokenizer kann mit einem Scanner im Kompilierungsprozess verglichen werden. Token sind in diesem Fall die Endpunkte von kontextfreien Grammatiken - Context-free grammar (CFG) - zur Analyse von Programmiersprachen-Terminalen.

Die Tokenisierung ist der erste Schritt in der NLP-Pipeline und kann daher den Rest der Pipeline stark beeinflussen. Der Tokenizer zerlegt unstrukturierte Daten, d. h. natürlichsprachliche Texte, in Informationseinheiten, die als einzelne Elemente gezählt werden können. Die so gezählte Anzahl von Token-Vorkommen in einem Dokument kann direkt als Vektor verwendet werden, der dieses Dokument repräsentiert. Ein solcher Ansatz ermöglicht es, aus einer unstrukturierten Zeichenfolge (einem Textdokument) unmittelbar eine für das maschinelle Lernen geeignete numerische Datenstruktur zu gewinnen. Diese Werte können den Computer direkt dazu veranlassen, nützliche Aktionen durchzuführen und Reaktionen zu erzeugen.

Die einfachste Art, einen Satz zu tokenisieren, ist die Verwendung von Leerzeichen als Worttrenner in Zeichenketten. Dies ist jedoch nicht optimal, denn wenn ein Satz z. B. ein Satzzeichen enthält, wird es von einem der Tokenisierer erfasst. Optimierte Tokenizer sind in mehreren Python-Bibliotheken implementiert, die jeweils ihre eigenen Vor- und Nachteile haben: spaCy, Stanford CoreNLP und Natural Language Toolkit (NLTK). NLTK und StanfordCoreNLP sind die am längsten bestehenden und am häufigsten verwendeten Bibliotheken zum Vergleich von NLP-Algorithmen in wissenschaftlichen Artikeln. Obwohl die StanfordCoreNLP-Bibliothek eine Python-Application Programming Interface (API) hat, basiert sie auf dem Java 8 CoreNLP-Anwendungsteil, der separat installiert und konfiguriert werden muss. Daher wurde in dieser

Arbeit der NLTK-Tokenizer verwendet.

Ein wichtiges Konzept im Tokenisierungsprozess sind N-Gramme. Ein N-Gramm ist eine Sequenz mit bis zu n Elementen, die aus einer Sequenz dieser Elemente, in der Regel einer Zeichenkette, extrahiert wurden. Im Allgemeinen können die Elemente eines N-Gramms Buchstaben, Silben, Wörter oder sogar Symbole sein. N-Gramme sind notwendig, weil bei der Konvertierung einer Menge von Wörtern in einen Vektor eine Folge von Token einen Großteil der Bedeutung verliert, die in der Reihenfolge dieser Wörter verkapselt ist. Wenn das Token-Konzept auf Mehrwort-Token, N-Gramme, ausgedehnt wird, kann die NLP-Pipeline einen erheblichen Teil der Bedeutung, die in der Wortfolge dieser Äußerungen enthalten ist, beibehalten. So bleibt beispielsweise das Wort „kein“, das die Bedeutung umkehrt, neben den benachbarten Wörtern stehen, wo es hingehört. Ohne N-Gramm-Tokenisierung würde ein solches Wort an verschiedenen Positionen herumhängen und seine Bedeutung würde mit dem gesamten Satz oder Dokument assoziiert werden, anstatt mit benachbarten Wörtern. Das Bigramm „war nicht“ behält viel mehr Bedeutung der einzelnen Wörter „war“ und „nicht“ als die entsprechenden Singlegramme im Multigrammvektor. Durch die Verknüpfung eines Wortes mit seinen Nachbarn in einem Förderband kann ein Teil seines Kontexts erhalten bleiben. N-Gramme sind also ein Instrument zur Speicherung von Kontextinformationen, während die Daten die Pipeline durchlaufen.

Die Größe des Vokabulars spielt eine wichtige Rolle für die Leistung der NLP-Pipeline. Die Größe des Wörterbuchs bestimmt die Größe der Trainingsstichprobe, die benötigt wird, um eine Überanpassung an ein bestimmtes Wort oder eine bestimmte Phrase zu vermeiden und die Größe der Trainingsmenge bestimmt die Kosten der Verarbeitung. Eine Technik zur Verringerung der Größe des Wörterbuchs besteht darin, Token, die ähnliche Dinge bedeuten, in einer einzigen normalisierten Form zusammenzufassen. Eine solche Technik reduziert die Anzahl der gespeicherten Token und verbessert die Verbindungen zwischen den Bedeutungen von Phrasen mit unterschiedlicher „Schreibweise“ der Tokens und verringert die Wahrscheinlichkeit des Überlernens. Eine der Normalisierungsmöglichkeiten ist die Groß- und Kleinschreibung - die Kombination mehrerer Schreibweisen eines Wortes, die sich nur in der Groß- und Kleinschreibung unterscheiden. In diesem Fall wird die Groß-/Kleinschreibung ignoriert und zwei identische Wörter, von denen eines groß und das andere klein geschrieben wird, werden als dasselbe Token behandelt.

Ein weiterer wichtiger Schritt bei der Tokenisierung von Texten ist die Entfernung von Stoppwörtern, um den Umfang des Wörterbuchs zu verringern. Stoppwörter sind in jeder Sprache gebräuchliche Wörter, die sehr häufig vorkommen, aber sehr viel weniger aussagekräftige Informationen über die Bedeutung eines Satzes enthalten (z. B. a, an, the, this, of, on im Englischen oder der, die, das, diese im Deutschen). Stoppwörter können jedoch nützliche Informationen enthalten, so dass man sie nicht immer verwerfen sollte. Das NLTK-Paket für Python enthält

derzeit die umfassendste Liste kanonischer Stoppwörter in verschiedenen Sprachen.

Eine weitere gängige Methode der Normalisierung ist die Beseitigung kleiner semantischer Unterschiede im Zusammenhang mit Pluralendungen und Possessivendungen von Wörtern oder sogar unterschiedlichen Verbformen. Diese Methode der Normalisierung, bei der ein gemeinsamer Wortstamm für verschiedene Wortformen gefunden wird, wird als Stemming bezeichnet. Der gemeinsame Wortstamm von „Gehäuse“ und „Haus“ ist zum Beispiel „Haus“. Beim Stemming werden Suffixe von Wörtern entfernt, um Wörter mit ähnlicher Bedeutung unter einem gemeinsamen Stamm zu gruppieren. Der Wortstamm muss nicht unbedingt ein gültiges Wort sein, sondern kann auch nur ein Token oder eine Bezeichnung sein, das mehrere mögliche Schreibweisen repräsentiert. Einer der Hauptvorteile des Stemming besteht darin, die Anzahl der Wörter zu verringern, deren Bedeutung das Sprachmodell im Auge behalten muss. Durch diese Methode wird die Größe des Wörterbuchs reduziert, wodurch der Verlust nützlicher Informationen und Bedeutungen so weit wie möglich begrenzt wird. Das Stemming spielt eine wichtige Rolle bei der Suche nach Schlüsselwörtern oder Informationen. Zwei der bekanntesten Algorithmen sind Porter's Stemmer und Snowball.

Mit Informationen über die Beziehungen zwischen den Bedeutungen verschiedener Wörter ist es möglich, mehrere Wörter miteinander zu verknüpfen, auch wenn ihre Schreibweise sehr unterschiedlich ist. Eine solche erweiterte Normalisierung eines Wortes auf seine semantische Wurzel - ein Lemma - wird Lemmatisierung genannt. Die Lemmatisierung ist potenziell eine viel genauere Art der Normalisierung als das Stemming oder die Groß- und Kleinschreibung, da sie die Bedeutung des Wortes berücksichtigt. Der Lemmatisierer verwendet eine Wissensbasis von Synonymen und Wortendungen, um nur eng verwandte Wörter zu einem Token zu kombinieren. In Python kann die Lemmatisierung mit dem NLTK-Paket implementiert werden, das den *WordNetLemmatizer* enthält.

Wie bereits gezeigt wurde, ist die Tokenisierung der Prozess der Zerlegung von Text in einzelne Wörter oder Token. Für viele Anwendungen der Verarbeitung natürlicher Sprache ist jedoch nicht nur wichtig, welche Wörter im Text enthalten sind, sondern man muss auch in der Lage sein, diese Wörter als Zahlen darzustellen. Dadurch wird es möglich, maschinelles Lernen und andere Algorithmen, die mit Zahlen arbeiten, zur Textverarbeitung einzusetzen. Bei der Verarbeitung natürlicher Sprache wird dazu die Wortvektorisierung verwendet.

2.3.3 Vektorisierung von Wörtern

Die Vektorisierung hilft dabei, eine numerische Darstellung zu finden, die die Bedeutung oder den Informationsgehalt der dargestellten Wörter widerspiegelt. Es gibt drei Möglichkeiten, Wörter und ihre Bedeutungen darzustellen:

- Multisets von Wörtern - Vektoren von Zahlen oder Häufigkeiten von Wörtern.
- Multisets von N-Grammen - Vektoren von Mengen von Wortpaaren (Bigramme), Worttripeln (Trigramme), usw.
- TF-IDF-Vektoren sind Indikatoren für Wörter, die deren Bedeutung am Besten darstellen.

Jede dieser Methoden kann entweder allein oder als Teil einer NLP-Pipeline angewendet werden. Alle beschriebenen Modelle sind statistisch in dem Sinne, dass sie auf Worthäufigkeiten beruhen.

Multisets von Wörtern

Bag Of Words (BOW)-Vektoren sind Vektoren, die Textdokumente als Wort-Multisets (Bags) darstellen, wobei ein Multiset eine Sammlung von Elementen ist, in der jedes Element mehr als einmal wiederholt werden kann. Ein Wort-Bag ist eine Darstellung von Text, die das Vorkommen von Wörtern in einem Dokument beschreibt. Sie beinhaltet zwei Dinge:

1. Ein Wörterbuch mit bekannten Wörtern.
2. Ein Maß für das Vorhandensein von bekannten Wörtern.

Dies wird als „Word-Bag“ bezeichnet, weil alle Informationen über die Reihenfolge oder Struktur der Wörter im Dokument verworfen werden. Das Modell interessiert sich nur dafür, ob bekannte Wörter im Dokument vorkommen, nicht aber, wo im Dokument. Bei diesem Ansatz wird das Histogramm der Wörter im Text betrachtet, d. h. wir betrachten jedes Wort als ein Merkmal [GOLDBERG 2017].

Ein Ansatz zur Vektorisierung einer Vielzahl von Wörtern besteht darin, für jedes Wort einen Einheitsvektor zu erstellen und diese Vektoren dann zu einem einzigen Vektor zu kombinieren. Ein unitärer Vektor für ein Wort ist ein Vektor, dessen Länge der Dimension des Raumes entspricht und alle seine Komponenten sind Null, mit Ausnahme einer Komponente, die der Position des Wortes im Wörterbuch entspricht. Unitäre Vektoren sind extrem spärlich. Enthält das Wörterbuch beispielsweise 1000 Wörter, dann hat jeder unitäre Vektor eine Dimension von 1000 und die Komponente, die der Position des Wortes im Wörterbuch entspricht, ist gleich eins. Zum Beispiel sieht der unitäre Vektor für den Satz „Mojito ist ein leckerer Cocktail“ (nach dem Tokenisierungsprozess) so aus:

```
Tokens: 'Cocktail ein ist leckerer Mojito'
unitärer Vektor: [[0, 0, 0, 0, 1]
                  [0, 0, 1, 0, 0]
                  [0, 1, 0, 0, 0]
                  [0, 0, 0, 1, 0]
                  [1, 0, 0, 0, 0]]
```

Abbildung 2.3: Unitärer Vektor

Eine der Eigenschaften der vektoriellen Darstellung von Wörtern und der tabellarischen Darstellung von Dokumenten ist, dass keine ursprünglichen Informationen verloren gehen. Das Originaldokument kann aus dieser einheitlichen Vektortabelle rekonstruiert werden. Daher werden solche unitären Vektoren in neuronalen Netzen, bei der sequence-to-sequence (seq2seq)-Modellierung und bei der Erstellung von Sprachmodellen verwendet. Unitäre Vektoren eignen sich hervorragend für alle NLP-Modelle oder -Pipelines, bei denen die Bedeutung des Ausgangstextes vollständig erhalten bleiben muss. Die Erstellung von unitären Vektoren für jedes Wort kann mit verschiedenen Bibliotheken in Python erfolgen. Zum Beispiel kann die NumPy-Bibliothek verwendet werden, um unitäre Vektoren zu erstellen.

Diese Vektordarstellung hat jedoch ihre Nachteile. Bei Dokumenten, die eine große Anzahl von Wörtern enthalten, wäre die Größe des unitären Vektorwörterbuchs enorm und würde enorme Ressourcen zur Verarbeitung erfordern. Außerdem macht es keinen Sinn, so viele Nullen zu speichern. Um mit dem Wörterbuch arbeiten zu können, ohne viel Zeit für die Verarbeitung zu verschwenden, ist es notwendig, das Wörterbuch zu verkleinern. Dazu ist es notwendig, auf die Möglichkeit der Rekonstruktion des Originaldokuments zu verzichten, d.h. sich die Wortfolge aller Dokumente zu merken. Es ist notwendig, den größten Teil der Informationen zu erfassen, nicht das gesamte Dokument. Dafür sind Multiset-Vektoren hilfreich.

Um einen Vektor für ein Multiset von Wörtern zu erstellen, müssen für jedes Wort der Menge ein Einheitsvektor erstellt und diese Vektoren dann zu einem einzigen Vektor kombiniert werden. Vektoren können nach verschiedenen Regeln kombiniert werden, z. B. durch Addition, Multiplikation, Verkettung oder andere Operationen. Wenn alle diese unitären Vektoren zusammengezählt werden, wird ein Vektor der Worthäufigkeit erhalten. Er spiegelt die Häufigkeit des Vorkommens von Wörtern wider, nicht die Reihenfolge, in der sie auftreten. Dieser Vektor kann verwendet werden, um ein ganzes Dokument oder einen Satz als einen einzigen Vektor von nicht sehr großer Länge darzustellen. Seine Länge entspricht der Länge des Wörterbuchs (der Anzahl der verfolgten eindeutigen Token). Außerdem kann bei einer einfachen Schlüsselwortsuche ein binärer Multiset-Vektor durch eine logische OR-Operation aus unitären Vektoren gewonnen werden.

Die Anzahl der Vorkommen eines Wortes in einem bestimmten Dokument wird als Termfrequenz (TF) bezeichnet. Je häufiger ein Wort in einem Dokument vorkommt, desto wahrscheinlicher ist es, dass das Dokument von diesem Thema handelt. Anstelle von reinen Wortzählungen können auch normalisierte Termhäufigkeiten zur Beschreibung von Dokumenten aus einem Korpus verwendet werden. Unter normalisierten Dokumenthäufigkeiten versteht man die Anzahl der Wörter im Verhältnis zur Länge des Dokuments. Die Normalisierung ist sehr nützlich, um die Bedeutung eines Dokuments zu ermitteln, da reine Häufigkeiten nicht immer die tatsächliche Situation widerspiegeln. Wenn beispielsweise das Wort „Cocktail“ 15 Mal in Dokument №1 und 400 Mal in Dokument №2 vorkommt, kann man daraus schließen, dass es in Dokument Nummer zwei eher um Cocktails geht. Vielleicht handelt es sich aber bei Dokument №1 um ein Rezept für die Zubereitung von Mai Tais mit 100 Wörtern, während Dokument №2 ein Roman mit 500.000 Wörtern ist, in dem die Figuren gelegentlich Cocktails trinken. Die Häufigkeit des Wortes „Cocktail“ in Dokument №1 beträgt also 0,15 und 0,0008 in Dokument №2, was die tatsächliche Bedeutung des Begriffs „Cocktail“ in diesen Dokumenten viel besser widerspiegelt:

$$\begin{aligned} \text{TF}(\text{Cocktail}, \text{Dokument}_1) &= 15/100 = 0,15 \\ \text{TF}(\text{Cocktail}, \text{Dokument}_2) &= 400/500000 = 0,0008 \end{aligned}$$

Abbildung 2.4: Normalisierte TF-Vektoren

Ebenso ist es möglich, diesen Wert für jedes Wort zu berechnen und die relative Bedeutung dieses Begriffs für jedes Dokument zu ermitteln, woraus sich die Bedeutung jedes Dokuments ableiten lässt. Um die Anzahl der Vorkommen von Wörtern zu ermitteln, kann das Objekt *collections.Counter* in Python verwendet und anschließend eine Normalisierung durchgeführt werden:

```
sentence = "Die mit der Getränkmischmaschine zubereiteten Cocktails  
können sich von denen an der Bar unterscheiden, da die Cocktails an  
der Bar stark mit Wasser verdünnt sind"  
tokenizer = TreebankWordTokenizer()  
tokens = tokenizer.tokenize(sentence.lower())  
token_counts = Counter(tokens)  
document_vector = []  
doc_length = len(tokens_intro)  
for key, value in token_counts.most_common():  
    document_vector.append(value/doc_length)  
>>> document_vector:  
[0.11538461538461539,  
 0.07692307692307693,  
 0.07692307692307693,  
 0.07692307692307693,  
 0.07692307692307693,  
 0.07692307692307693,  
 0.038461538461538464,  
 ...,  
 0.038461538461538464]
```

Abbildung 2.5: Normalisierte TF-Vektoren mit Python

Wenn mehr als ein Dokument bearbeitet wird, wird für jedes Dokument ein eigener Häufigkeitsvektor erstellt, der jedoch absolut alle Begriffe aus allen Dokumenten (alle Wörter aus dem Wörterbuch) enthält, auch wenn einige der Wörter nicht in allen Dokumenten vorkommen. In diesem Fall werden die Vektoren an diesen Stellen Nullwerte enthalten.

Das ultimative Ziel der Textvektorisierung besteht darin, natürliche Sprache in ein Format zu übersetzen, das der Computer verstehen kann. Eine Möglichkeit, dies zu erreichen, ist die Verwendung von Vektorräumen. Vektorräume sind mathematische Objekte, in denen jedes Element durch einen Vektor dargestellt wird und verschiedene Operationen mit Vektoren ermöglichen es uns, nützliche Merkmale zu extrahieren. So wird ein Vektor mit zwei Werten in einem zweidimensionalen Vektorraum verortet, ein Vektor mit drei Werten in einem dreidimensionalen Vektorraum und so weiter. Beispielsweise ist der Vektorraum von Dokumenten, die nur aus den Wörtern „Cocktail“ und „Maschine“ bestehen, zweidimensional und die Vektoren der enthaltenen Dokumente können wie folgt dargestellt werden:

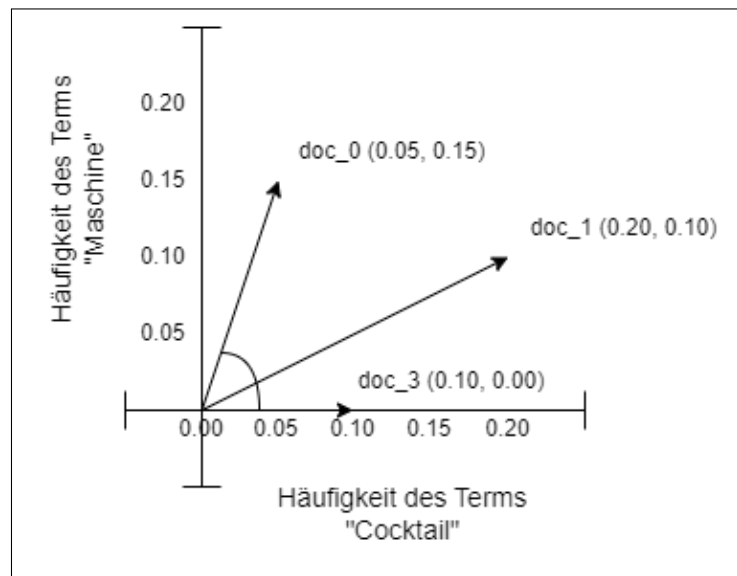


Abbildung 2.6: Zweidimensionale Vektoren von Termhäufigkeiten

Für einen Vektorraum von Dokumenten in einer natürlichen Sprache ist die Dimensionalität des Vektorraums gleich der Anzahl der verschiedenen Wörter im gesamten Korpus. Jedes Dokument in einem K-dimensionalen Vektorraum kann durch einen K-dimensionalen Vektor beschrieben werden. Wenn es Vektorrepräsentationen für alle Dokumente im gleichen Raum gibt, dann ist es möglich, sie zu vergleichen. Der euklidische Abstand zwischen Vektoren kann gemessen werden, indem sie subtrahiert werden und die Länge des Abstands zwischen ihnen berechnet wird, den sogenannten Abstand im Sinne der L2-Norm (L2-Metrik). Dies ist der Abstand auf einer geraden Linie von dem Ort, der durch die Spitze (Scheitelpunkt) des einen Vektors gegeben ist, zu dem Ort, der der Spitze des anderen Vektors entspricht.

$$d(p, q) = \|q - p\|_2 = \sqrt{(q_1 - p_1)^2 + \dots + (q_n - p_n)^2} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

Abbildung 2.7: Euklidische Norm

Der Ochiai-Koeffizient (Kosinuskoeffizient) ist ein Maß für die Ähnlichkeit zwischen zwei Multisets. Der Ochiai-Koeffizient kann verwendet werden, um den Grad der Ähnlichkeit zwischen zwei Dokumenten auf der Grundlage ihrer Wort-Multisets zu bestimmen. Er ist einfach der Kosinus des Winkels zwischen zwei Vektoren (Theta).

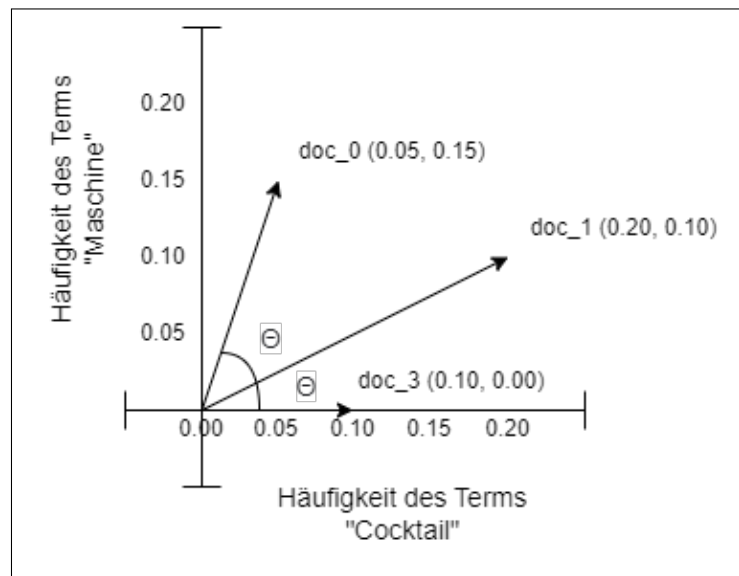


Abbildung 2.8: Zweidimensionale Theta

Er kann mithilfe des euklidischen Skalarprodukts nach folgender Formel berechnet werden:

$$\mathbf{A} \cdot \mathbf{B} = \|\mathbf{A}\| \|\mathbf{B}\| \cos \theta$$

Abbildung 2.9: Euklidisches Skalarprodukt

$$\text{cosine similarity} = S_C(A, B) := \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}},$$

Abbildung 2.10: Ochiai-Koeffizient

Das Skalarprodukt zweier Vektoren kann also berechnet werden, indem die Elemente der Vektoren paarweise multipliziert und dann die Ergebnisse addiert werden. Anschließend wird das Ergebnis durch die Norm jedes Vektors dividiert. Die Norm eines Vektors ist gleich dem euklidischen Abstand zwischen seinem Kopf und seinem Ende - der Quadratwurzel aus der Summe der Quadrate seiner Elemente. Dieses normierte Skalarprodukt liegt wie der Kosinuswert zwischen -1 und 1. Es ist auch der Kosinus des Winkels zwischen den beiden Vektoren. Dieser Wert ist gleich dem Anteil des längeren Vektors, der von der senkrecht dazu stehenden Projektion des kürzeren Vektors abgedeckt wird. Daraus ist ersichtlich, wie weit unsere beiden Vektoren in dieselbe Richtung zeigen.

Ein offensichtlicher Vorteil ist, dass der vom Ochiai-Koeffizienten angenommene Wertebereich

(von -1 bis +1) für die meisten Aufgaben des maschinellen Lernens geeignet ist. Ein Kosinuskoeffizient von 1 spiegelt identische normalisierte Vektoren wider, die über alle Dimensionen hinweg in eine ähnliche Richtung zeigen. Die Längen dieser Vektoren können unterschiedlich sein, aber sie zeigen in dieselbe Richtung. Je näher der Kosinuskoeffizient bei 1 liegt, desto kleiner ist der Winkel zwischen den Vektoren. Im NLP enthalten Dokumente, deren Kosinuskoeffizient der Vektoren nahe bei 1 liegt, ähnliche Wörter im gleichen Verhältnis. Daher ist es wahrscheinlicher, dass es sich bei Dokumenten, deren Vektoren nahe beieinander liegen, um dieselbe Sache handelt. Wenn der Kosinuskoeffizient 0 ist, haben die Vektoren keine gemeinsamen Komponenten. Sie stehen in allen Dimensionen in einem 90° -Winkel zueinander. Bei TF-NLP-Vektoren tritt diese Situation nur ein, wenn die beiden Dokumente keine gemeinsamen Wörter haben. Da sie völlig unterschiedliche Wörter verwenden, sprechen sie von völlig unterschiedlichen Dingen. Das bedeutet nicht unbedingt, dass die Bedeutung oder das Thema unterschiedlich ist, sondern nur, dass sie unterschiedliche Wörter verwenden. Bei einem Wert des Kosinuskoeffizient von -1 sind die Vektoren völlig entgegengesetzt. Sie zeigen in entgegengesetzte Richtungen. Bei einfachen Vektoren der Termhäufigkeiten ist dies jedoch nicht möglich, da Vektoren der Termhäufigkeiten nicht negativ sein können. Sie befinden sich also immer im gleichen Quadranten des Vektorraums.

Neben den Multisets von Wörtern gibt es weitere Vektorisierungsmethoden, die ebenfalls in der Verarbeitung natürlicher Sprache verwendet werden - Multisets von N-Grammen und TF-IDF-Vektoren. Bei Multisets von N-Grammen handelt es sich um Mengen von mehreren aufeinanderfolgenden Wörtern in einem Text. N-Gramme können helfen, den Kontext und die Beziehung zwischen Wörtern in einem Text zu erfassen. In diesem Fall wird nicht die Häufigkeit der einzelnen Begriffe im Dokument berücksichtigt, sondern die Häufigkeit des Auftretens von N-Grammen. TF-IDF-Vektoren sind jedoch eine noch effizientere Methode zur Vektorisierung von Text. Im Gegensatz zu Multisets von Wörtern und N-Grammen berücksichtigen TF-IDF-Vektoren nicht nur die Häufigkeit von Wörtern in einem Dokument, sondern auch die Bedeutung jedes Wortes im Kontext des gesamten Textkorpus. So haben Wörter, die in allen Dokumenten häufig vorkommen, eine geringere Bedeutung als seltene Wörter, die dazu beitragen, ein Dokument von anderen zu unterscheiden.

TF-IDF-Vektoren

Das TF-IDF-Vektorkonzept ist eines der grundlegenden Werkzeuge der natürlichen Sprachverarbeitung. Es wird verwendet, um Textdokumente in numerische Vektoren umzuwandeln, die für Textanalyse, Dokumentenklassifikation, Empfehlungssysteme und viele andere Anwendungen verwendet werden können.

Eine der wichtigsten Fragen, die sich bei der Verarbeitung natürlicher Sprache stellen, ist die Bestimmung der Bedeutung von Wörtern in einem Text. Die Bedeutung von Wörtern kann

anhand ihrer Häufigkeit in einem Dokument bestimmt werden. Einige Wörter kommen jedoch häufig in allen Dokumenten vor, so dass sie nicht allzu wichtig sein sollten. Andererseits können einige seltene Wörter von großer Bedeutung sein, weil sie möglicherweise nur in einem bestimmten Dokument vorkommen. TF-IDF-Vektoren lösen dieses Problem, indem sie nicht nur die Häufigkeit der Wörter in einem Dokument, sondern auch ihre Bedeutung im Kontext des gesamten Textkorpus berücksichtigen. Für jedes Wort in einem Dokument kann sein TF-IDF-Wert nach der folgenden Formel berechnet werden:

$$TF(t, d) = \frac{Anzahl(t)}{Anzahl(d)} \quad (2.1)$$

$$IDF(t, D) = \log \frac{Anzahl(D)}{Anzahl(D, \text{ die } t \text{ beinhalten})} \quad (2.2)$$

$$TFIDF(t, d, D) = TF(t, d) \cdot IDF(t, D) \quad (2.3)$$

Dabei ist t ein Wort, d ein Dokument, $TF(t, d)$ der Wert der Worthäufigkeit im Dokument und $IDF(t, D)$ der Wert der inversen Dokumenthäufigkeit. Um einen TF-IDF-Vektor für jedes Dokument zu erhalten, müssen die TF-IDF-Werte für alle Wörter, die in dem Dokument vorkommen, verwendet werden. Dies ergibt einen Vektor, bei dem jede Komponente einem TF-IDF-Wert für ein bestimmtes Wort im Dokument entspricht. Beispielsweise werden die folgenden drei Dokumente betrachtet:

Dokument №1 „Ein Cocktail ist variationsreich und auch zuhause gut zu mixen“

Dokument №2 „Der Mojito ist ein klassischer Cocktail und ein erfrischender Drink der immer passt“

Dokument №3 „Dank unserem Rezept machen Sie diesen Cocktail selber zuhause“

Zunächst müssen die Dokumente tokenisiert und unnötige Stoppwörter herausgefiltert werden, die für das Verständnis der Bedeutung der Dokumente keine große Rolle spielen. Nach diesem Prozess sehen die Dokumente wie folgt aus:

Dokument №1 „Cocktail variationsreich zuhause gut mixen“

Dokument №2 „Mojito klassischer Cocktail erfrischender Drink immer passt“

Dokument №3 „Dank Rezepte machen Cocktail selber zuhause“

Als erstes müssen die normalisierten TF-Vektoren jeden Terms berechnet werden:

TF(t, d)	Dokument №1	Dokument №2	Dokument №3
Cocktail	$1/5 = 0.2$	$1/7 = 0.143$	$1/6 = 0.167$
variationsreich	0.2	0	0
zuhaus	0.2	0	0.167
gut	0.2	0	0
mixen	0.2	0	0
Mojito	0	0.143	0
klassischer	0	0.143	0
erfrischender	0	0.143	0
Drink	0	0.143	0
immer	0	0.143	0
passt	0	0.143	0
Dank	0	0	0.167
Rezepte	0	0	0.167
machen	0	0	0.167
selber	0	0	0.167

Tabelle 2.1: TF-Vektoren

Anschließend können die IDF-Vektoren für jeden Term (über den gesamten Dokumentenkörper) berechnet werden:

IDF(t, D)	Dokumentenkörper
Cocktail	$\log(3/3) = 0$
variationsreich	0.477
zuhaus	0.176
gut	0.477
mixen	0.477
Mojito	0.477
klassischer	0.477
erfrischender	0.477
Drink	0.477
immer	0.477
passt	0.477
Dank	0.477
Rezepte	0.477
machen	0.477
selber	0.477

Tabelle 2.2: IDF-Vektoren

Jetzt kann mit der Berechnung der TF-IDF-Vektoren fortgefahren werden:

TF-IDF(t, d, D)	Dokument №1	Dokument №2	Dokument №3
Cocktail	$0.2 \cdot 0 = 0$	$0.143 \cdot 0 = 0$	$0.167 \cdot 0 = 0$
variationsreich	0.0954	0	0
zuhaus	0.0352	0	0.03
gut	0.0954	0	0
mixen	0.0954	0	0
Mojito	0	0.068	0
klassischer	0	0.068	0
erfrischender	0	0.068	0
Drink	0	0.068	0
immer	0	0.068	0
passt	0	0.068	0
Dank	0	0	0.08
Rezepte	0	0	0.08
machen	0	0	0.08
selber	0	0	0.08

Tabelle 2.3: TF-IDF-Vektoren

Wie aus dem Beispiel ersichtlich, haben Wörter, die in allen Dokumenten vorkommen, wie z.B. „Cocktail“, in jedem Dokument den gleichen TF-IDF-Wert, Wörter, die in mehr als einem Dokument vorkommen, wie z.B. „zuhaus“, haben einen niedrigeren IDF-Wert als seltene Wörter, wie z.B. „selber“ oder „Rezept“, die in den Dokumenten, in denen sie vorkommen, einen hohen IDF-Wert und damit einen hohen TF-IDF-Wert haben. Diese durch die Berechnung von TF-IDF erhaltene Zahl assoziiert also ein bestimmtes Wort oder Token mit einem bestimmten Dokument in einem bestimmten Korpus und legt dann einen numerischen Wert für die Bedeutung dieses Wortes in diesem Dokument fest, wobei sein Vorkommen im gesamten Korpus berücksichtigt wird. Da für jedes Wort in den Dokumenten nun TF-IDF-Werte vorliegen, können diese zu TF-IDF-Vektoren für jedes Dokument kombiniert werden. Der TF-IDF-Vektor für jedes Dokument hat eine Dimension, die der Gesamtzahl der eindeutigen Wörter im Textkorpus entspricht. Jedes Element im Vektor entspricht einem bestimmten Wort und der Wert des Elements ist gleich seinem TF-IDF-Wert für dieses Dokument. In diesem Beispiel sehen die TF-IDF-Vektoren für jedes Dokument wie folgt aus:

Dokument №1 [0. 0.0954 0.0352 0.0954 0.0954 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]

Dokument №2 [0. 0. 0. 0. 0. 0.068 0.068 0.068 0.068 0.068 0.068 0. 0. 0. 0.]

Dokument №3 [0. 0. 0.03 0. 0. 0. 0. 0. 0. 0. 0. 0.08 0.08 0.08 0.08]

Aus diesen Vektoren kann die Ähnlichkeit zwischen den Dokumenten berechnet werden, z. B. mit Hilfe des Kosinuskoeffizienten. Je näher die Werte des Kosinuskoeffizienten bei 1 liegen, desto ähnlicher sind sich die Dokumente. Für die maximale derzeit erreichbare Effizienz bei der semantischen Suche, der Dokumentenklassifikation, bei Dialogsystemen und den meisten anderen Anwendungen ist es daher ausreichend, nur einfache TF-IDF-Vektoren in die Pipeline einzugeben. TF-IDF-Vektoren sind die erste Stufe der Pipeline, die Grundmenge der aus dem Text extrahierten Merkmale. Auf der Grundlage der TF-IDF-Vektoren können Themenvektoren berechnet werden, die für die semantische Analyse von Texten erforderlich sind.

Es gibt viele Bibliotheken, um die TF-IDF-Vektorisierung in Python zu implementieren. Einige der beliebtesten Bibliotheken sind Scikit-Learn, Gensim, NLTK und SpaCy.

2.3.4 Syntaktische Analyse

Mit Hilfe von TF-IDF-Vektoren ist es möglich, die Bedeutung von Wörtern in einem Textfragment zu schätzen. TF-IDF-Vektoren und -Matrizen können verwendet werden, um zu verstehen, wie wichtig jedes Wort für die Gesamtbedeutung eines Textfragments in einer Reihe von Dokumenten ist. Solche TF-IDF-Bedeutungsschätzungen sind nicht nur für Wörter, sondern auch für kurze Wortfolgen, N-Gramme, möglich. Und letztere sind sehr praktisch für die Suche im Text (wenn die gesuchten Wörter oder N-Gramme genau bekannt sind).

Experimente auf dem Gebiet des NLP haben zu einem Algorithmus geführt, mit dem die Bedeutung von Wortkombinationen ermittelt und Vektoren zur Darstellung dieser Bedeutung berechnet werden können. Dieser Algorithmus wird latente semantische Analyse (Latent Semantic Analysis (LSA)) genannt. Mit ihr lassen sich nicht nur Wörter, sondern ganze Dokumente als Vektoren darstellen [RYGL u. a. 2017]. Diese Vektoren werden semantische Vektoren oder Themenvektoren genannt. Diese Themenvektoren können verwendet werden, um nach Dokumenten anhand ihrer semantischen Bedeutung zu suchen - eine semantische Suche. In den meisten Fällen sind die Ergebnisse dieser Suche wesentlich besser als die einer Schlagwortsuche (TF-IDF-Suche). Manchmal liefert eine semantische Suche genau die Dokumente, nach denen der Benutzer gesucht hat, auch wenn er nicht die richtigen Wörter für seine Anfrage gefunden hat. Darüber hinaus können diese semantischen Vektoren verwendet werden, um die Wörter und N-Gramme zu finden, die das Thema einer Aussage, eines Dokuments oder eines Korpus (Menge von Dokumenten) am besten repräsentieren. Mit Hilfe der Wortvektoren und ihrer relativen Bedeutung ist es möglich, die repräsentativsten Wörter für ein Dokument zu finden, wie z. B. eine Reihe von Schlüsselwörtern, die seine Bedeutung widerspiegeln. Es ist nun auch möglich, zwei beliebige Aussagen/Dokumente zu vergleichen und festzustellen, wie nahe sie sich in ihrer Bedeutung sind.

Themenvektoren sind eine Möglichkeit, ein Dokument als Vektor darzustellen, wobei jedes

Element einem Thema entspricht. Jedes Thema ist eine Wahrscheinlichkeitsverteilung auf die Wörter im Korpus und für jedes Dokument kann geschätzt werden, wie sehr es mit jedem Thema zu tun hat. Zur Berechnung der Themenvektoren wird der LSA-Algorithmus verwendet, der auf einer SWZ (eng. Singular Value Decomposition (SVD)) der Term-Dokumenten-Matrix beruht.

Bei der SWZ wird die Matrix in drei Komponenten zerlegt: die Matrix der linken singulären Vektoren, die Matrix der rechten singulären Vektoren und die Diagonalmatrix der singulären Werte. Diese drei einfacheren Matrizen aus der SWZ offenbaren jedoch Eigenschaften der ursprünglichen TF-IDF-Matrix. Sie eignen sich zur Vereinfachung der Matrix. Man kann die Matrizen abschneiden (einige Zeilen und Spalten weglassen), bevor man sie erneut multipliziert, wodurch die Anzahl der Dimensionen, mit denen man sich im Vektorraummodell befassen muss, reduziert wird. Das Ergebnis der Multiplikation dieser abgeschnittenen Matrizen ist nicht dieselbe TF-IDF-Matrix, sondern eine verbesserte: die neue Dokumentendarstellung enthält das Wesentliche, die latente Semantik dieser Dokumente. Aus diesem Grund wird SWZ auch in anderen Anwendungen wie der Komprimierung eingesetzt (z. B. zur Komprimierung von JPEG-Bildern).

Die latente semantische Analyse (LSA) ist eine mathematische Methode, um den besten Weg zur linearen Transformation (Rotation und Streckung) eines beliebigen Satzes von NLP-Vektoren, wie TF-IDF-Vektoren oder Multiset-Vektoren, zu finden. Für viele Anwendungen ist es am besten, die Achsen (Dimensionen) in den neuen Vektoren an der größten Streuung (Varianz) der Worthäufigkeiten auszurichten [JURAFSKY u. a. 2009]. Es können dann diejenigen Dimensionen aus dem neuen Vektorraum ausgeschlossen werden, die nur geringe Auswirkungen auf die Vektoränderungen von Dokument zu Dokument haben. Diese Anwendung der SWZ wird als trunkierte SWZ (trunkated SVD) bezeichnet.

LSA verwendet SWZ, um nach Wortkombinationen zu suchen, die gemeinsam für die größte Varianz in den Datenwerten verantwortlich sind. Die TF-IDF-Vektoren können so gedreht werden, dass ihre neuen Dimensionen (Basisvektoren) mit diesen Richtungen der maximalen Varianz übereinstimmen. Die Basisvektoren stellen die Koordinatenachsen des neuen Vektorraums dar. Jede der Dimensionen (Koordinatenachsen) wird zu einer Kombination von Worthäufigkeiten und nicht zu einer einzelnen Worthäufigkeit. Sie können also als gewichtete Kombinationen der Wörter betrachtet werden, die die verschiedenen im Korpus verwendeten Themen bilden. Die LSA-Methode liefert, wie der IDF-Teil von TF-IDF, Informationen darüber, welche Vektordimensionen für die Semantik (Bedeutung) von Dokumenten eine wichtige Rolle spielen. Die Dimensionen (Topics) mit der geringsten Variation in den Vektoren von Dokument zu Dokument können verworfen werden. Solche Themen mit geringer Varianz stellen in der Regel Rauschen dar und sind eine Ablenkung für jeden Algorithmus des maschinellen Lernens. Wenn ein Thema

in allen Dokumenten etwa gleich häufig vorkommt und nicht zur Unterscheidung zwischen den Dokumenten beiträgt, kann es verworfen werden, so dass die Vektordarstellung verallgemeinert werden kann, so dass sie sich besser für Dokumente eignet, auf die die Pipeline noch nicht gestoßen ist, selbst wenn sie aus einem anderen Kontext stammen.

SWZ ist ein Algorithmus zur Zerlegung einer beliebigen Matrix in drei Faktoren, d. h. drei Matrizen, die multipliziert werden, um die ursprüngliche Matrix zu bilden. Dieser Algorithmus wird in der linearen Algebra verwendet, um eine Matrix umzukehren und die drei mit SWZ berechneten Matrixmultiplikatoren haben mehrere nützliche mathematische Eigenschaften, die für die Dimensionalitätsreduktion und die LSA-Methode geeignet sind [ECKART und YOUNG 1936].

Jedes Mal, wenn SWZ eine Term-Dokumenten-Matrix BOW oder eine Term-Dokumenten-Matrix TF-IDF durchläuft, findet SWZ übereinstimmende Wortkombinationen. Um diese übereinstimmenden Wörter zu finden, berechnet SWZ die Korrelation zwischen den Spalten (Termen) der Term-Dokumenten-Matrix. SWZ ermittelt gleichzeitig die Korrelation der Verwendung von Begriffen in Dokumenten und die Korrelation von Dokumenten untereinander. Auf der Grundlage dieser Informationen berechnet SWZ auch lineare Kombinationen von Begriffen mit maximaler Varianz innerhalb eines Korpus. Solche Kombinationen von Termhäufigkeiten werden zu Themen. Und nur die Themen, die die maximale Information und Varianz innerhalb eines Korpus beherbergen, werden als relevant markiert und beibehalten, die anderen werden abgeschnitten. Darüber hinaus besteht die Möglichkeit der linearen Transformation (Drehung) der Term-Dokument-Vektoren, um sie in kürzere Themenvektoren für jedes der Dokumente umzuwandeln. [WONG, ZIARKO und WONG 1985]

SWZ gruppiert die Begriffe mit hoher Korrelation (weil sie häufig in denselben Dokumenten vorkommen) und auch gleichzeitig in einer Reihe von Dokumenten variieren [ISBELL und VIOLA 1998]. Diese linearen Wortkombinationen werden als Themen betrachtet, die BOW-Vektoren (oder TF-IDF-Vektoren) in Themenvektoren verwandeln, die angeben, auf welche Themen sich das Dokument bezieht. Ein Themenvektor ist eine Zusammenfassung (Synthese) des Inhalts des Dokuments. Eine SWZ sieht wie folgt aus [DEERWESTER u. a. 1990]:

$$A_{m \times n} = U_{m \times p} S_{p \times p} V_{p \times n}^T \quad (2.4)$$

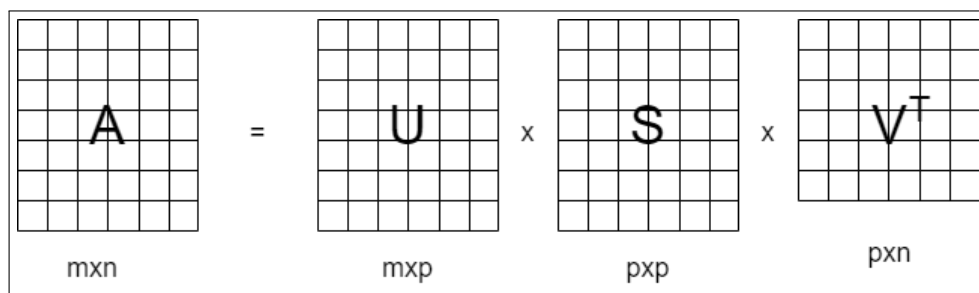


Abbildung 2.11: SWZ-Matrize

Wobei m die Anzahl der Begriffe im Wörterbuch, n die Anzahl der Dokumente im Korpus und p die Summe der Themen im Korpus ist, die der Anzahl der Wörter entspricht.

Nach dem Theorem von Eckart-Young reduziert die SWZ das Rauschen und die Spärlichkeit der ursprünglichen Matrix, indem sie durch eine Matrix derselben Dimensionalität, aber niedrigeren Rangs ersetzt wird, in der nur die wichtigsten Informationen erhalten bleiben. [ECKART und YOUNG 1936]

Die U-Matrix enthält die Term-Themen-Matrix. Sie wird auch als linke singuläre Vektormatrix bezeichnet, da sie die Zeilenvektoren enthält, die von links mit der Spaltenvektormatrix zu multiplizieren sind. Die U-Matrix spiegelt die gegenseitige Korrelation von Wörtern und Themen auf der Grundlage des gemeinsamen Auftretens von Wörtern im selben Dokument wider. Vor der Trunkierung (Entfernung der Spalten) ist sie quadratisch. Die Anzahl der Spalten und Zeilen in ihr entspricht der Anzahl der Wörter im Wörterbuch (m). Vor der Trunkierung ist die Anzahl der Themen (p) gleich der Anzahl der Wörter. Die U-Matrix enthält als Spalten Vektoren der Themen für alle Wörter des Korpus. Das bedeutet, dass sie verwendet werden kann, um einen Wort-Dokument-Vektor (TF-IDF-Vektor oder BOW-Vektor) in einen Themen-Dokument-Vektor umzuwandeln. Um einen neuen „topic-document“-Vektor zu erhalten, multipliziert man einfach die U-Matrix „topic-word“ mit einem beliebigen „word-document“-Vektor. Der Punkt ist, dass die Gewichte (Scores) in den Zellen der U-Matrix die Bedeutung der Wörter für die Themen widerspiegeln.

Die quadratische diagonale S-Matrix (Sigma-Matrix) enthält die Singulärwerte (Eigenwerte) des Themas. Diese singulären Werte geben an, wie viel Information in jeder Dimension des neuen semantischen (Themen-)Vektorraums enthalten ist. In der diagonalen Matrix befinden sich Werte, die nicht Null sind, nur auf der Diagonalen von der linken oberen Ecke bis zur rechten unteren Ecke. Alle anderen Positionen in der S-Matrix enthalten Nullen. Und die Dimensionen (Themen) sind so angeordnet, dass die erste Dimension die meisten Informationen (erklärbare Varianz) über den Fall enthält. Wenn es also notwendig ist, das Themenmodell abzuschnei-

den, kann man mit den Messungen von unten rechts beginnen und sich nach links bewegen. Und dann aufhören, diese singulären Werte auf Null zu setzen, wenn der Fehler des thematischen Modells eine signifikante Auswirkung auf den Gesamtfehler der NLP-Pipeline haben wird.

Die V^T -Matrix „Dokument-Dokument“ enthält als Spalten die rechten singulären Vektoren und gibt Aufschluss über die gemeinsamen Bedeutungen von Dokumenten, da sie die Häufigkeit der Verwendung der gleichen Themen in Dokumenten in dem neuen semantischen Modell von Dokumenten misst. Es enthält so viele Zeilen (p) und Spalten, wie es Dokumente im Korpus gibt.

Nachdem die singuläre Zerlegung abgeschlossen ist, müssen die Themen so gekürzt werden, dass sie die Bedeutung der Dokumente wiedergeben. Dabei müssen die Elemente der Matrix S in absteigender Reihenfolge sortiert werden. In diesem Fall kann die Formel in einer anderen Form geschrieben werden:

$$A_k = U_k S_k V_k^T \quad (2.5)$$

Dabei sind U_k und V_k Matrizen, die durch Extraktion der ersten k Spalten aus den Matrizen U bzw. V gewonnen werden.

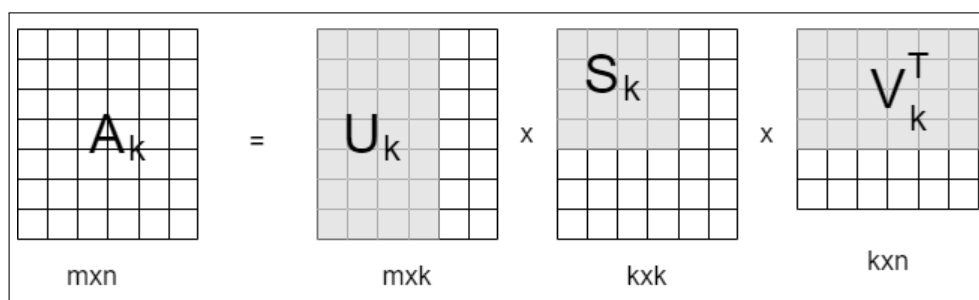


Abbildung 2.12: Trunkierte SWZ-Matrize

Im latent-semantischen Modell kann also jeder Text und jeder Begriff durch Vektoren in einem reduzierten Raum der Dimension k - dem Themenraum - dargestellt werden. Diese SWZ wird als parsimonious bezeichnet, weil sie in dem Fall, in dem k viel kleiner als m, n und p ist, eine erhebliche Kompression der ursprünglichen Information ermöglicht. Komprimierung ist in dem Sinne zu verstehen, dass ein Teil der von der ursprünglichen Matrix übermittelten Informationen verloren geht und nur die wichtigsten (dominanten) Informationen erhalten bleiben. Der Informationsverlust entsteht durch die Vernachlässigung kleiner Singulärzahlen. Je mehr Singulärzahlen verworfen werden, d. h. je kleiner k, desto größer ist dieser Verlust. Mit anderen Worten: Der Informationsverlust ist umso größer, je niedriger der Rang der approximierenden Matrix ist. In diesem Fall gehen zunächst die nicht signifikanten Werte verloren, wodurch die dominanten Werte der Matrix erhöht werden.

Der SWZ-Algorithmus, das Herzstück von LSA, erkennt, ob Wörter immer nebeneinander vorkommen und fasst sie im Betreff zusammen. Auf diese Weise kann er die Anzahl der Messungen ohne zusätzlichen Aufwand reduzieren. LSA (SWZ) ist eine großartige Methode, um Wort-Dokument-Matrizen zu komprimieren und potenzielle zusammengesetzte Wörter oder N-Gramme für die Pipeline zu erkennen. LSA reduziert die Anzahl der Messungen und hilft so, Übertraining zu vermeiden. Sie führt eine Verallgemeinerung auf der Grundlage eines kleinen Datensatzes durch, indem sie eine lineare Beziehung zwischen den Worthäufigkeiten annimmt.

LSAs können auch für die semantische Suche verwendet werden. Eine semantische Suche (semantic search) ist eine Art der Volltextsuche, die die Bedeutung der Wörter in der Anfrage und der gesuchten Dokumente berücksichtigt. Durch die Verwendung von Themenvektoren ist es möglich, die Bedeutung von Wörtern, Dokumenten, Aussagen und Korpora zu vergleichen. Es ist möglich, Cluster von ähnlichen Dokumenten und Aussagen zu finden. Mit Themenvektoren wird der Abstand zwischen Dokumenten nicht nur auf der Grundlage der in ihnen verwendeten Wörter betrachtet, sondern auf der Grundlage von Wörtern, die der Bedeutung entsprechen. Semantische Suchen mit Themenvektoren beschränken sich nicht nur auf die Suche nach Schlüsselwörtern und Relevanzrankings, die allein auf der Auswahl von Wörtern oder Vokabeln basieren. Sie können verwendet werden, um Dokumente zu finden, die für die Anfrage relevant sind, und nicht nur solche, die gut zu den Wortstatistiken passen. Dazu muss die Suchanfrage des Nutzers zunächst mit denselben Attributen vektorisiert werden wie die Dokumentenvektorisierung. Dann kann SWZ verwendet werden, um die Dimensionalität des Abfragevektors auf die gleichen Dimensionen wie die der Dokumentvektoren zu reduzieren. Anschließend kann der Kosinuskoeffizient zwischen dem Abfragevektor und den Dokumentvektoren berechnet werden, um festzustellen, welche Dokumente der Abfrage am nächsten liegen. LSA und SWZ können also zum Aufspüren versteckter Themen in Textdokumenten sowie zur semantischen Suche in großen Dokumentensammlungen verwendet werden.

Es gibt viele Python-Bibliotheken und -Werkzeuge für die Implementierung von LSA und SWZ. Eine der beliebtesten Bibliotheken für Matrix- und Singulärzerlegung ist NumPy. NumPy erleichtert die Durchführung von Matrixoperationen und verfügt über eine Funktion zur Singulärzerlegung: `numpy.linalg.svd()`. Zusammen mit anderen Bibliotheken wie *Pandas* und *Scikit-learn* kann LSA effizient durchgeführt werden. Um LSA zu starten, müssen die Textdaten mit Hilfe des TF-IDF-Modells vektorisiert werden, wie zuvor beschrieben. Dann muss eine singuläre Zerlegung der Term-Dokument-Matrix mit `numpy.linalg.svd()` durchgeführt werden. Danach kann eine bestimmte Anzahl von Hauptkomponenten (Topics) ausgewählt werden, die zur Darstellung der Dokumente verwendet werden. Dies kann durch die Auswahl der Komponenten mit den höchsten Singulärwerten geschehen. Anhand dieser Komponenten kann dann die Nähe zwischen

den Dokumenten bestimmt werden.

Scikit-learn bietet auch eine *TruncatedSVD*-Klasse, mit der eine Truncated-SVD durchgeführt werden kann, um nur die ausgewählten Hauptkomponenten zu erhalten. Diese Klasse erleichtert die Durchführung von LSA auf großen Dokumentenkorpora und verfügt über zusätzliche Funktionen wie Datenverarbeitung und Normalisierung sowie automatische Auswahl der Anzahl der Komponenten auf der Grundlage des Prozentsatzes der erklärten Varianz.

2.3.5 Ansätze für die Erstellung eines Chatbots

Derzeit gibt es vier Hauptansätze für die Erstellung eines Chatbots [LANE, HOWARD und HAPKE 2019]:

- Musterabgleich: Musterabgleich und Antwortvorlagen (vorgefertigte Antworten)
- Grounding: logische Wissensgraphen und das Ziehen von Schlussfolgerungen aus diesen basierend auf diesen Graphen
- Suche: Abrufen von Text
- Generierungsmethoden: Statistik und maschinelles Lernen

Die vier grundlegenden Ansätze zur Erstellung von Chatbots lassen sich kombinieren, was zu benutzerfreundlicheren Chatbots führt. Eine Vielzahl von Anwendungen nutzen alle vier grundlegenden Methoden. Hybride Chatbots unterscheiden sich hauptsächlich darin, wie genau sie diese Ansätze kombinieren und wie viel Gewicht auf jeden einzelnen Ansatz gelegt wird.

Musterabgleich

Bei den ersten Chatbots basierte die Antwort auf die Nachricht eines Benutzers auf einem Mustervergleich. Diese Chatbots suchen nach Mustern im eingehenden Text und geben eine feste (gemusterte) Antwort, wenn eine Übereinstimmung gefunden wird [WOUDENBERG 2014].

Solche rudimentären Dialogsysteme sind vor allem in automatisierten Benutzerunterstützungssystemen mit interaktiven Sprachmenüs nützlich, wo es möglich ist, das Gespräch an einen Menschen weiterzuleiten, wenn der Chatbot keine Antwortmuster mehr hat.

Da es viele NLP-Dienstprogramme in Python-Paketen gibt, ist es möglich, komplexere Chatbots auf der Grundlage von Mustervergleichen zu erstellen, indem man die Bot-Logik nach und nach direkt in Python mit regulären Ausdrücken und Suchmustern aufbaut.

1995 machte sich Richard Wallace daran, einen allgemeinen Rahmen für die Erstellung von

Chatbots auf der Grundlage des Pattern-Matching-Ansatzes zu schaffen. Zwischen 1995 und 2002 schuf seine Entwicklergemeinschaft die AIML zur Beschreibung von Mustern und Chatbot-Antworten.

AIML ist eine deklarative Sprache, die auf dem eXtensible Markup Language (XML)-Standard basiert, der die Sprachkonstrukte und Datenstrukturen einschränkt, die im Bot verwendet werden dürfen. [AIML Foundation o. D.] Ein Chatbot, der auf AIML basiert, sieht folgendermaßen aus:

```
<?xml version="1.0" encoding="UTF-8"?><aiml version="2.0">
  <category>
    <pattern>Hi, Bot</pattern>
    <template>Hallo</template>
  </category>
  <category>
    <pattern>Wie geht es dir, bot?</pattern>
    <template>Gut</template>
  </category>
</aiml>
```

Abbildung 2.13: AIML Chatbot

Eine der Einschränkungen von AIML ist die Art der Muster, die abgeglichen werden können und auf die reagiert wird. Der AIML-Kern (Pattern Matching Engine) reagiert nur auf Eingabetext, der einem vom Entwickler manuell vorgegebenen Muster entspricht. Unschärfe Suchanfragen, Smileys, Satzzeichen, Tippfehler oder falsch geschriebene Wörter sind nicht erlaubt, es findet kein automatischer Abgleich statt. In AIML müssen alle Synonyme manuell einzeln beschrieben werden.

Grounding

Die Grounding-Methode ist ein Ansatz zur Erstellung eines Chatbots auf der Grundlage logischer Wissensgraphen und der Durchführung von Schlussfolgerungen auf der Grundlage dieser Graphen. [DIANA und ISMAEL 2011] Sie wird verwendet, um natürliche Sprache zu verarbeiten und sie dem Verständnis des Bots zuzuordnen. Das Wesentliche an der Grounding-Methode ist, dass der Chatbot nicht nur die Textnachrichten, sondern auch den Kontext und die Umgebung verarbeitet, um Anfragen besser zu verstehen und zu beantworten. Durch die Extraktion von Informationen wird ein Netz von Verbindungen oder Fakten geschaffen. Dieses Netz logischer Verbindungen

zwischen Entitäten - ein Graph oder eine Wissensbasis - kann die Grundlage für die Antworten des Chatbots bilden.

Ein Beispiel für eine Grounding-Methode ist die Verwendung eines Wissensgraphen zur Beschreibung der Umgebung. Ein Wissensgraph enthält Informationen über die Objekte, mit denen der Bot interagieren kann und die Beziehungen zwischen ihnen. Ein Wissensgraph könnte zum Beispiel Informationen über ein Glas auf einem Tisch und das darin befindliche Wasser enthalten. Wenn ein Benutzer eine Frage stellt, verwendet der Chatbot den Wissensgraphen, um den Kontext der Anfrage zu verstehen und die am besten geeignete Antwort abzuleiten. Wenn ein Benutzer zum Beispiel fragt: „Wie hoch ist die Temperatur des Wassers in dem Glas auf dem Tisch?“, kann der Chatbot Informationen aus dem Wissensgraphen verwenden, um die Frage zu beantworten.

Ein solcher Wissensgraph kann abgeleitet werden, um Fragen über die in dieser Wissensbasis enthaltene Welt zu beantworten und anschließend können auf der Grundlage der logischen Antworten die Werte der in den Antworten enthaltenen Template-Variablen ausgefüllt werden, um natürlichsprachliche Antworten zu erstellen. Ursprünglich wurden auf diese Weise Systeme zur Beantwortung von Fragen eingerichtet, wie z. B. der Watson-Bot von IBM (heutzutage wird für ähnliche Systeme jedoch die Informationssuchemethode verwendet). Der Wissensgraph stellt eine Art „Erdung“ des Chatbots in der realen Welt dar.

Die Erstellung von Chatbots auf der Grundlage von „Grounding“ eignet sich hervorragend für Chatbots, die Fragen generieren, bei denen das zur Beantwortung einer Frage erforderliche Wissen in einer umfangreichen Wissensbasis enthalten ist, die aus einer offenen Datenbank (z. B. Wikidata, Open Mind Common Sense oder DBpedia) bezogen werden kann.

Einer der Hauptvorteile der Grounding-Methode besteht darin, dass sie sich an ein sich veränderndes Umfeld anpassen kann. Wenn der Benutzer zum Beispiel ein Glas Wasser von einem Tisch auf einen anderen stellt, wird der Wissensgraph automatisch aktualisiert, um diese Änderung widerzuspiegeln.

Die Grounding-Methode hat jedoch auch ihre Grenzen. So kann es vorkommen, dass bei der Verarbeitung großer Informationsmengen Zusammenhänge nicht berücksichtigt werden und dem Bot möglicherweise verborgen bleiben.

Insgesamt ist die Grounding-Methode ein effektiver Ansatz zur Erstellung wissensbasierter Chatbots. Sie ermöglicht es dem Bot, Benutzeranfragen besser zu verstehen und eine genauere Antwort zu geben.

Suche

Die Informationssuchemethode ist eine der Methoden zum Aufbau von Chatbots, die auf der Extraktion von Informationen aus einer großen Menge von Textinformationen basiert. Die Hauptidee der Informationssuchemethode ist die Analyse des Eingabetextes (Benutzeranfrage), die Auswahl von Schlüsselwörtern und Phrasen daraus und die anschließende Suche nach den relevantesten Informationen in der Wissensdatenbank oder in offenen Quellen. [DIANA und ISMAEL 2011]

Die Wissensbasis kann auch eine Art „Gesprächsprotokoll“ sein, in Form von Aussage-Antwort-Paaren. Dabei sucht der Bot nach früheren Aussagen in den Protokollen früherer Unterhaltungen. Der Bot kann nicht nur in den Protokollen seiner eigenen Gespräche suchen, sondern auch in beliebigen Transkripten von Gesprächen zwischen Menschen, Gesprächen zwischen Menschen und Bots oder sogar Gesprächen zwischen Bots. Aber wie immer gilt: je besser die Eingabedaten, desto besser das Ergebnis. Daher ist es notwendig, die Datenbank früherer Gespräche sorgfältig zu säubern und zu organisieren, damit der Bot nach einem qualitativ hochwertigen Gespräch sucht und es dann imitiert.

Für die Umsetzung der Informationssuchemethode werden verschiedene Algorithmen und Techniken verwendet, z. B. Indizierung und Schlagwortsuche, Kontextsuche, Textanalyse mit Hilfe von maschinellen Lernverfahren usw. Die Informationssuchemethode kann in Python mit verschiedenen Bibliotheken und Tools wie NLTK, Scikit-learn und Gensim implementiert werden.

Einer der ersten Schritte bei der Implementierung einer Informationssuchemethode in Python ist die Vorbereitung der Daten. Dies erfordert Tokenisierung, Lemmatisierung und die Entfernung von Stopp-Wörtern. Als nächstes muss ein Index auf der Grundlage von Schlüsselwörtern erstellt werden. Der Index kann auf der Grundlage von Bag-of-Words oder TF und IDF (TF-IDF-Modelle) erstellt werden. Sobald der Index erstellt ist, kann eine Stichwortsuche durchgeführt werden. Dazu muss die Benutzeranfrage in einen Vektor umgewandelt und mit den Dokumentvektoren im Index verglichen werden. Dies kann mit Hilfe der Scikit-learn-Bibliothek erfolgen [SCIKIT-LEARN o. D.] Sobald die relevantesten Dokumente gefunden wurden, können sie in eine Rangfolge gebracht und als Antwort auf die Benutzeranfrage angezeigt werden.

Der Vorteil der Informationssuchemethode besteht darin, dass sie ein schnelles und genaues Auffinden der gewünschten Informationen ermöglicht, insbesondere wenn die Wissensbasis gut strukturiert ist und genügend Informationen enthält. Ein Nachteil dieser Methode ist jedoch, dass sie den Kontext der Anfrage nicht berücksichtigt und nicht immer eine vollständige und genaue Antwort auf die Frage des Nutzers liefert. Wenn die Aussage semantisch mit der vom Bot zu beantwortenden übereinstimmt, ist es möglich, die Antwort wortwörtlich und ohne Änderungen wiederzuverwenden. Aber selbst wenn die Datenbank alle möglichen Benutzeräußerungen enthält,

wird der Bot die Persönlichkeiten der Personen widerspiegeln, die diese Äußerungen machen. Wenn die Antworten konsistent sind und von einer Vielzahl von Personen stammen, ist das gut. Problematisch wird es jedoch, wenn die Äußerung, auf die der Bot reagieren soll, vom Gesamtkontext des jeweiligen Gesprächs oder von den Umständen in der Umgebung abhängt, die sich seit der Erstellung des Dialogkorpus geändert haben können.

Beispielsweise sollte der Bot auf die Frage „Wie spät ist es?“ nicht die von der Person gegebene Antwort, sondern die am besten geeignete Aussage aus der Datenbank verwenden. Diese Antwort funktioniert nur, wenn die Zeit, zu der die Frage gestellt wurde, mit der Zeit übereinstimmt, zu der die passende Äußerung aus der Datenbank aufgezeichnet wurde. Neben dem natürlichsprachlichen Text der Äußerung müssen auch ähnliche Informationen über die Zeit - der Kontext (Zustand) - erfasst und verglichen werden. Sie spielt vor allem dann eine wichtige Rolle, wenn die Semantik der Äußerung auf eine aktive Veränderung des im Kontext (Wissensbasis des Chatbots) erfassten Zustands hinweist.

Um den Zustand (Kontext) in einem Chatbot auf der Grundlage der Informationssuche zu berücksichtigen, kann etwas Ähnliches für einen Chatbot mit Musterabgleich durchgeführt werden, da die Auflistung einer Liste von Benutzeraussagen nur eine andere Art ist, ein Muster zu beschreiben. Dies auch ist der Ansatz von Amazon Lex [AMAZON o. D.] und Google Dialogflow [CHAWLA o. D.] Anstatt ein starres Muster zu beschreiben, um den Befehl des Benutzers zu erfassen, können der Dialogflow-Engine einfach ein paar Beispiele geliefert werden. So wie jedes Muster im Chatbot auf der Grundlage der Musterzuordnung einem Zustand zugeordnet wurde, muss auch hier nur die Aussage-Antwort-Beispielpaare mit dem genannten Zustand verknüpft werden.

Der suchbasierte Chatbot indiziert also den Korpus der Dialoge, so dass er leicht frühere Aussagen finden kann, die derjenigen ähnlich sind, auf die er antworten muss und antwortet dann mit einer der passenden Aussagen aus dem Korpus, die er sich „gemerkt“ und für eine schnelle Suche indiziert hat. Im Allgemeinen ist die Methode der Informationssuche eine der gängigsten und beliebtesten Methoden zum Aufbau von Chatbots, die in verschiedenen Bereichen wie Wirtschaft, Medizin, Tourismus und vielen anderen eingesetzt werden.

Generierungsmethoden

Generierungsmethoden sind einer der wichtigsten Ansätze bei der Entwicklung von Chatbots auf der Grundlage künstlicher Intelligenz. Sie ermöglichen es Chatbots, Textantworten auf der Grundlage der Analyse der eingehenden Nachricht und des Kontextes des Dialogs zu generieren. Die folgenden Generierungsmodelle sind nützlich, um einen kreativen Chatbot zu erstellen, der Dinge sagen kann, die noch niemand zuvor gesagt hat:

- Sequenz-zu-Sequenz-Konvertierungsmodelle: Modelle, die darauf trainiert sind, Antworten auf der Grundlage von Eingabesequenzen zu generieren;
- Restricted Boltzmann Machines (RBM): Markov-Ketten, die so trainiert werden, dass sie die „Energie“-Funktion minimieren [NUPUR SHARMA o. D.];
- Generative Adversarial Network (GAN): statistische Modelle, die darauf trainiert sind, einen Experten, der die Qualität eines Gesprächs bewertet, zu täuschen. [LI u. a. 2017]

Die Vorteile des Einsatzes der Generierungsmethoden:

- Flexibilität: Generative Methoden können für eine breite Palette von Aufgaben eingesetzt werden, einschließlich Texterstellung, Sprachübersetzung, Verarbeitung natürlicher Sprache und mehr.
- Automatisierung: Generative Methoden können auf großen Datensätzen trainiert werden, wodurch die Erstellung von Inhalten automatisiert werden kann.
- Qualität: Generative Methoden zeigen eine hohe Qualität bei der Textgenerierung, Sprachübersetzung und anderen Aufgaben der natürlichen Sprachverarbeitung, wenn sie auf einem ausreichend großen Datensatz trainiert werden.
- Schnelligkeit: Generative Methoden können schneller arbeiten als Menschen, was die Erstellung von Inhalten mit großer Geschwindigkeit ermöglicht.

Die Nachteile der generativen Methoden:

- Große Datenmengen für das Training: Generative Methoden benötigen große Datenmengen für das Training, was bei einigen Aufgaben schwierig sein kann, insbesondere wenn nur ein kleiner Datensatz zur Verfügung steht.
- Sicherheitsrisiken: Generative Methoden können Inhalte erzeugen, die möglicherweise falsch, unvollständig oder irreführend sind. Dies kann zu Sicherheitsrisiken führen, wenn der generierte Inhalt für wichtige Entscheidungen verwendet wird.
- Unterstützungsbedarf: Generative Methoden können erhebliche Unterstützung benötigen, um effektiv zu sein. Dies kann die Modellabstimmung, die Auswahl optimaler Parameter und die Optimierung der Modellleistung auf einer bestimmten Hardwarekonfiguration umfassen.
- Modellbeschränkungen: Generative Methoden können Beschränkungen hinsichtlich der Arten von Inhalten haben, die sie erzeugen können, insbesondere wenn sie nur auf bestimmte Datentypen trainiert wurden.

Eine der beliebtesten Methoden zur Texterstellung ist die sequence-to-sequence-Methode (seq2seq). Die seq2seq-Methode basiert auf Recurrent Neural Networks (RNN), die die Simulation von Datenfolgen ermöglichen. Sie besteht aus zwei Hauptteilen: einem Encoder und einem Decoder. Ein Encoder empfängt eine Wortfolge und baut daraus einen Kontextvektor auf, der Informationen über die Eingabedaten enthält. Der Decoder erhält diesen Vektor als Eingabe und beginnt mit der Generierung einer Folge von Antwortnachrichten, wobei er schrittweise den Kontext und die zuvor generierten Wörter berücksichtigt. [ALAMMAR o. D.]

Einer der Hauptvorteile der seq2seq-Methode ist ihre Fähigkeit, qualitative und grammatikalisch korrekte Textantworten zu generieren, einschließlich Antworten, die nicht in den Trainingsdaten enthalten waren. Sie kann auch mit langen Sequenzen umgehen, was sie ideal für die Generierung von Antworten in Dialogsystemen macht. Darüber hinaus kann die seq2seq-Methode in einer Vielzahl von Anwendungen eingesetzt werden, z. B. in der maschinellen Übersetzung, der Spracherkennung und anderen.

Die seq2seq-Methode hat jedoch ihre Nachteile. Sie erfordert große Datenmengen zum Trainieren und Verarbeiten sowie erhebliche Rechenressourcen. Dies kann die Anwendung der Methode bei einigen Anwendungen einschränken. Wenn der Trainingsdatensatz nicht eine ausreichend große Bandbreite möglicher Antworten repräsentiert, kann das Modell außerdem dazu neigen, vorhersehbare oder falsche Antworten zu erzeugen.

Die Implementierung der seq2seq-Methode in Python kann mit der TensorFlow-Bibliothek erfolgen, die eine Reihe von Werkzeugen für den Aufbau und das Training neuronaler Netze bietet. In TensorFlow kann man die vortrainierten seq2seq-Modelle verwenden oder ein eigenes Modell erstellen, indem die Architektur und die Trainingsparameter des Netzwerks konfiguriert wird. [TENSORFLOW o. D.[b]]

2.3.6 Mehrschichtiges Perzeptron für ein Mehrklassen-Klassifizierungsproblem

Ein Mehrschichtiges Perzeptron (MLP) ist eine neuronale Netzwerkarchitektur, die aus mehreren Schichten von Neuronen besteht und für Klassifizierungs- und Regressionsaufgaben verwendet werden kann. Mehrschichtige Perzeptronen sind Arten von neuronalen Netzen, die bidirektional sind, da sie die Eingabedaten vorverteilen und die Gewichte zurückverteilen. Ein mehrschichtiges Perzeptron hat eine Eingabeschicht und eine Ausgabeschicht mit einer oder mehreren versteckten Schichten. Bei MLP sind alle Neuronen einer Schicht mit allen Neuronen der nächsten Schicht verbunden. Die Eingabeschicht empfängt Eingangssignale und die gewünschte Aufgabe wird von der Ausgabeschicht ausgeführt. Die versteckten Schichten sind für alle Berechnungen zuständig. Die Architektur des mehrschichtigen Perzeptrons sieht folgendermaßen aus:

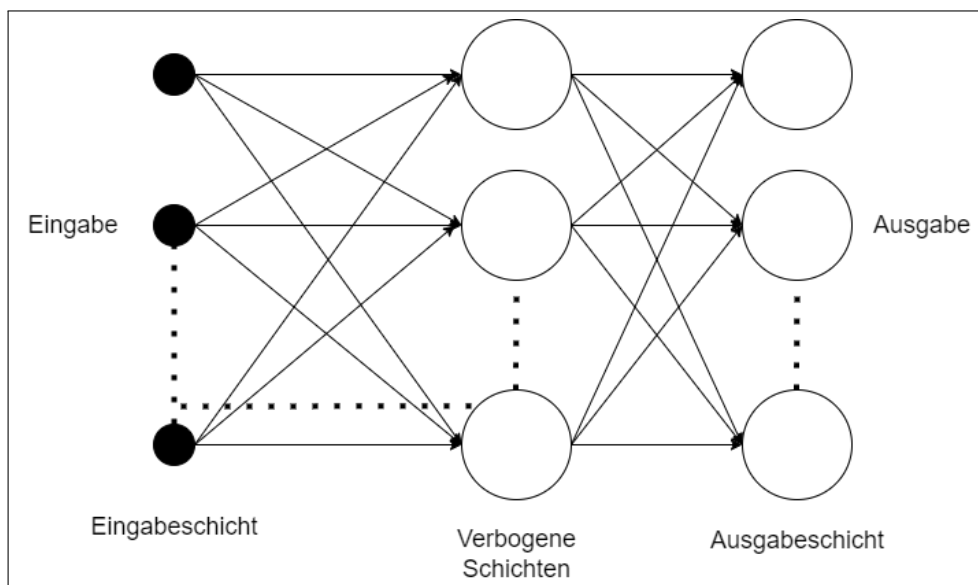


Abbildung 2.14: Mehrschichtiges Perzeptron

Für ein Mehrklassen-Klassifizierungsproblem kann ein MLP mit einer Verlustfunktion wie der kategorialen Kreuzentropie (Categorical Crossentropy) trainiert werden. Zunächst muss ein MLP über eine Eingabeschicht verfügen, die die Merkmale eines Objekts aufnimmt, und eine Ausgabeschicht, die die Wahrscheinlichkeiten für die Zugehörigkeit des Objekts zu den einzelnen Klassen ausgibt. Die Dimensionalität der Ausgabeschicht muss der Anzahl der zu klassifizierenden Klassen entsprechen. Außerdem können zwischen der Eingabe- und der Ausgabeschicht mehrere verborgene Schichten liegen, die nichtlineare Transformationen der Eingabedaten durchführen, so dass das MLP-Modell komplexe Beziehungen zwischen Objektmerkmalen und Klassen erfassen kann. Jedes Neuron im MLP nimmt als Eingabe die gewichtete Summe der Neuronenausgaben der vorherigen Schicht, fügt eine Vorspannung hinzu und durchläuft dann eine Aktivierungsfunktion. Die Aktivierungsfunktion wird verwendet, um dem MLP Nichtlinearität hinzuzufügen und kann z. B. sigmoid-Funktion oder ReLU sein. Beim MLP-Training mit Gradientenabstieg werden die Gewichte und Vorspannungen der Neuronen bei jeder Trainingsiteration aktualisiert, wobei der Wert der Verlustfunktion minimiert wird.

Die kategoriale Kreuzentropie ist eine Verlustfunktion, die für die Klassifizierung mehrerer Klassen verwendet wird und die Diskrepanz zwischen den wahren Klassenbezeichnungen und den vorhergesagten Klassenwahrscheinlichkeiten misst. Daher kann MLP nach dem Training zur Klassifizierung neuer Objekte verwendet werden, indem die Objektmerkmale durch MLP geleitet werden und die Ausgabewahrscheinlichkeiten für jede Klasse ermittelt werden und dann die Klasse mit der höchsten Wahrscheinlichkeit als vorhergesagte Klassenbezeichnung ausgewählt wird.

Die Theorie hinter der kategorialen Kreuzentropie ist die Informationstheorie und die Messung der Unsicherheit. Angenommen, es gibt ein System, das mit einem probabilistischen Modell beschreiben werden muss. Dieses System habe N mögliche Ergebnisse, die jeweils mit den Wahrscheinlichkeiten p_1, p_2, \dots, p_N auftreten können, wobei $\sum_i p_i = 1$ und $p_i \geq 0$ für $i = 1, 2, \dots, N$. Dann kann die Entropie H des Systems wie folgt definiert werden:

$$H = - \sum (p_i * \log_2 p_i) \quad (2.6)$$

Dabei wird $\sum (p_i * \log_2 p_i)$ über i von 1 bis N genommen [KOECH 2022]. Diese Formel zeigt, dass die Entropie des Systems ein Maß für die Ungewissheit oder Unsicherheit bei der Wahl eines der möglichen Ergebnisse ist.

Wenn die kategoriale Kreuzentropie als Verlustfunktion für ein Mehrklassen-Klassifizierungsproblem verwendet wird, wird erwartet, dass das trainierte Modell einen Ausgabewahrscheinlichkeitsvektor hat, bei dem jedes Element des Vektors einer der Klassen entspricht. Je unsicherer der Modellausgabevektor für ein bestimmtes Beispiel ist, desto höher ist der Wert der kategorialen Kreuzentropie. Umgekehrt ist der Wert der kategorialen Kreuzentropie gleich Null, wenn das Modell eine Klasse für ein bestimmtes Beispiel genau vorhersagt. Beim Training von Klassifizierungsmodellen für Aufgaben des maschinellen Lernens mit mehreren Klassen wird daher versucht, den Wert der kategorialen Kreuzentropie zwischen dem Ausgabevektor des Modells und der richtigen Klasse für jedes Eingabebeispiel zu minimieren. Zu diesem Zweck wird ein Optimierungsverfahren verwendet, z. B. den Gradientenabstieg, um die Gewichte und Verzerrungen des Modells bei jeder Trainingsiteration zu aktualisieren.

In Keras kann die kategoriale Kreuzentropie als Verlustfunktion für verschiedene Modelltypen wie mehrschichtige Perzeptrone, faltungsneuronale Netze und rekurrente neuronale Netze verwendet werden. Es handelt sich um eine Standardverlustfunktion für Mehrklassen-Klassifizierungsprobleme. Bei der Verwendung der kategorialen Kreuzentropie in Keras können auch verschiedene Parameter wie Klassengewichte und Modellbewertungsmetriken angepasst werden, um genauere Klassifizierungsergebnisse zu erhalten. Somit ist die kategoriale Kreuzentropie ein wichtiges Werkzeug beim maschinellen Lernen zur Lösung von Klassifizierungsproblemen mit mehreren Klassen. Ihre Verwendung ermöglicht es, Modelle zu trainieren, die Eingabebeispiele genau klassifizieren und eine hohe Genauigkeit bei Klassifizierungsaufgaben erreichen können.

Kapitel 3

Anforderungen

Im Folgenden sollen die Anforderungen an das Ergebnis der Arbeit konkretisiert werden.

3.1 Antwortzeit

Diese Eigenschaft beschreibt die Zeitdauer vom registrieren eines Sprachbefehls bis zur Ausführung des Befehls durch die Mischmaschine und das Zurückgeben einer Antwort an den Benutzer. Die Antwortzeit spielt eine große Rolle bei der Bedienbarkeit eines interaktiven Systems, um das es sich bei der Mischmaschine handelt. Zu lange Antwortzeiten können dazu führen, dass der Benutzer seine ursprünglichen Ziele vergisst oder in Stress gerät, da in den aller meisten Fällen der Grund für eine lange Antwortzeit vor dem Benutzer verborgen bleibt. Umgekehrt können zu kurze Antwortzeiten ebenfalls zu Stress und Fehlbedienung seitens des Benutzers führen. Dies liegt unter anderem daran, dass kurze Antwortzeiten den Benutzer dazu veranlassen weniger über seine Aktionen und deren Folgen nachzudenken. Als eine für viele Anwendungen geeignete Antwortzeit werden zwei bis vier Sekunden genannt [HERCZEG 2018]. Als maximale tolerierbare Antwortzeit werden für dieses Projekt sechs Sekunden festgelegt. Diese vergleichsweise lange Zeitdauer lässt sich zum Einen mit den langwierigen aber notwendigen Berechnungen begründen, die für die Spracherkennung und -verarbeitung benötigt werden. Zum Anderen werden die Auswirkungen einer langen Antwortzeit als gering eingeschätzt, da der Benutzer für diese Anwendung keine Teilarbeitsschritte o. ä. im Gedächtnis behalten muss. Das Ziel des Benutzers sich ein Getränk zubereiten zu lassen ist nach dem Eingang des Sprachbefehls bereits erfüllt.

3.2 Offline-Funktionalität

Die Mischmaschine sollte für die Sprachsteuerung keine Verbindung zum Internet benötigen, da dies die möglichen Einsatzorte der Maschine deutlich einschränken würde. Diese Anforderung schränkt die möglichen, einzusetzenden Technologien zur Spracherkennung und -verarbeitung stark ein, da keine Cloud-Services wie bspw. *Google Cloud Speech* eingesetzt werden können

[GOOGLE o.D.] Eine weitere Herausforderung die dadurch entsteht ist, dass Berechnungen die unter Umständen sehr aufwendig sein können nicht ausgelagert sondern auf der Hardware innerhalb der Mischmaschine ausgeführt werden müssen.

3.3 Lautstärke

Die Lautstärke der, von der Mischmaschine zurückgegebenen Antwort, muss laut genug sein, sodass sie vom Benutzer gut verstanden werden kann. Diese Eigenschaft schränkt die Art und Weise wie die Hardware (Computer und Mikrofon) in die Mischmaschine eingebaut werden kann ein und welche Art von Hardware überhaupt verwendet werden kann.

3.4 Entfernung

Mit dieser Eigenschaft ist die Entfernung des Anwenders zu der Mischmaschine gemeint. Es muss dem Anwender ermöglicht werden aus einer moderaten Entfernung mit der Mischmaschine über die Sprachsteuerung zu interagieren. Sowohl die Eingabe eines Befehls über die Sprachsteuerung als auch die zu hörende Antwort sollte mindestens aus einer Entfernung von einem Meter möglich sein. Dafür müssen die Lautsprecher eine bestimmte Lautstärke erreichen können und das Mikrofon eine moderate Empfindlichkeit aufweisen.

3.5 Antworten

Die Antworten, die durch die Mischmaschine an den Benutzer zurückgegeben werden, sollen mit Hilfe eines eigens erstellten Sprachmodells auf Basis von künstlicher Intelligenz und ML erfolgen. Die Antworten der Mischmaschine sollen außerdem bissiger bzw. sarkastischer Natur sein was, je nach verwandter Technik, bei der Auswahl der Trainingsdaten eine große Rolle spielt.

Des Weiteren bestehen die Anforderungen, dass der Benutzer auf Deutsch mit der Mischmaschine kommunizieren können muss und die Antworten der Mischmaschine kontrollierbar sein müssen. Mit der Kontrollierbarkeit ist gemeint, dass Vorhersagen darüber gemacht werden können, was die Mischmaschine in Etwa auf eine bestimmte Frage oder sonstige Benutzereingabe antworten wird. Dies soll verhindern, dass der Benutzer von unerwarteten Reaktionen seitens der Maschine überrascht wird und die Antworten der Maschine den Benutzer nicht beleidigen (aufgrund der sarkastischen bzw. humorvollen Art und Weise, wie die Maschine antworten soll).

3.6 Kosten

Die Materialkosten sollten einen gewissen Maximalbetrag nicht überschreiten. Zu den benötigten Materialien zählen ein Mikrocomputer zur Durchführung der Spracherkennung und -verarbeitung,

ein Mikrofon zur Aufnahme der Sprache und Lautsprecher zur Tonausgabe. Das Ziel besteht darin einen Betrag von 200€ nicht zu überschreiten.

3.7 Verbrauch von Arbeits- und Festplattenspeicher

Diese Anforderung korreliert mit der Anforderung nach Offline-Funktionalität. Diese bedingt, dass aufwendige, rechen- oder speicherintensive Operationen nicht auf entfernte Rechner ausgelagert werden können sondern, alles auf „kleiner“ Hardware innerhalb der Mischmaschine ablaufen muss. Diese Einschränkung soll anhand eines konkreten Beispiels verdeutlicht werden. Die vierte Version der bekannten Mikrocomputerreihe *Raspberry Pi* umfasst im Modell B maximal acht Gigabyte Read Only Memory (RAM) und einen Micro-SD-Karten-Steckplatz [LTD, RASPBERRY PI o. D.] Zwar sind Micro-SD-Karten von bis zu mehreren hundert Gigabyte erhältlich, jedoch muss dabei der Kostenfaktor mit beachtet werden. Dadurch sind sowohl Primär- als auch Sekundärspeicher stark beschränkt.

Auch Prozessorgeschwindigkeit und Grafikkartenleistung können bei der Spracherkennung und -verarbeitung eine Rolle spielen. Diese ist bei Mikrocomputern ebenfalls eingeschränkt und korreliert negativ mit der Anforderung an Geschwindigkeit. Einen ebenso negativen Einfluss auf die Geschwindigkeit des Gesamtsystems haben die virtuelle Vergrößerung des RAM durch Swapping oder die SD-Karte selbst, welche im Vergleich mit anderen Speichertechnologien in Sachen Geschwindigkeit deutlich das Nachsehen hat.

In anbetracht der geschilderten Herausforderungen und Kostenbetrachtung wird der maximal zu verbrauchende Arbeits- und Festplattenspeicher für dieses Projekt auf acht Gigabyte RAM und 20 Gigabyte Festplattenspeicher festgelegt.

3.8 Anpassungsfähigkeit

Diese Anforderung beschreibt den Grad der Einfachheit bei Änderung der Anforderungen oder Umwelt die Sprachsteuerung der Mischmaschine an diese neuen Begebenheiten anzupassen. Ein einfaches Beispiel für die Anpassungsfähigkeit des Systems wäre das Hinzukommen eines Behälters innerhalb der Mischmaschine. Tritt dieser Fall ein sollte es leicht möglich sein die Sprachsteuerung so anzupassen, dass der Benutzer auch die Möglichkeit bekommt aus dem neuen, fünften Behälter Getränke zu Mischen und zu Bestellen.

3.9 Bewertung der Ansätze zur Erstellung eines Dialogsystems

Unter Berücksichtigung der genannten Anforderungen und der im Kapitel 2.3.5 besprochenen Ansätze zur Erstellung eines Dialogsystems, können die folgenden Schlussfolgerungen gezogen werden, ob die Ansätze den Anforderungen entsprechen oder nicht:

Ansatz/ Anforderung	Musterabgleich	Grounding	Suche	Generierungsmethoden
Antwortzeit	entspricht	entspricht nicht	entspricht	entspricht nicht
Offline-Funktionalität	entspricht	entspricht	entspricht	entspricht
Sarkastische Antworten	entspricht	entspricht	entspricht	entspricht nicht
Antworten auf deutsch	entspricht	entspricht nicht	entspricht nicht	entspricht nicht
kontrollierbare Antworten	entspricht	entspricht	entspricht	entspricht nicht
Verbrauch von Arbeits- und Festplattenspeicher	entspricht	entspricht nicht	entspricht nicht	entspricht nicht

Tabelle 3.1: Bewertung der Ansätze zur Erstellung eines Dialogsystems

Kapitel 4

Konzept

In diesem Kapitel wird zunächst erläutert, wie die Steuerung der Getränkemischmaschine durch Sprachbefehle im Allgemeinen ablaufen wird. Anschließend werden mehrere Konzepte vorgestellt, die das allgemeine Konzept konkretisieren. Diese werden anhand der, in Kapitel 4.2 erläuterten Kriterien, bewertet. Zuletzt wird die Wahl des finalen Konzepts begründet.

4.1 Allgemein

Der Benutzer soll über Spracheingaben mit der Mischmaschine interagieren können. Dafür muss das Gesprochene zunächst durch ein Mikrofon aufgenommen werden. Anschließend können die Audiosignale weiterverarbeitet werden. Der Benutzer soll hierbei nicht auf fest vorgegebene Sprachbefehle beschränkt sein, sondern für nahezu jede Eingabe eine sinnvolle Antwort zurück erhalten. Um dies zu gewährleisten wird die Spracheingabe durch ein Sprachmodell, welches mittels maschinellen Lernverfahren trainiert wurde, verarbeitet. Ergebnisse dieser Verarbeitung sind die Antwort, die an den Benutzer zurückgegeben wird, und ein konkreter Befehl für die Mischmaschine. Ein Beispiel für einen solchen Befehl könnte etwa die Zubereitung eines bestimmten Getränks sein. Für die Ausgabe einer Antwort ist ein Lautsprecher notwendig. Denkbar wäre auch eine textbasierte Ausgabe, allerdings ginge damit der Eindruck des Benutzers verloren eine echte Konversation mit der Mischmaschine zu führen. Das Sprachmodell mit der Getränkemischmaschine zu verknüpfen stellt eine Herausforderung dieser Arbeit dar.

4.2 Bewertungskriterien

Im Folgenden sind die Bewertungskriterien für die einzelnen Konzepte aufgelistet:

- Freiheitsgrade in der Spracheingabe des Benutzers: Erhält der Benutzer passende Antworten zurück egal was er sagt oder ist er auf einige wenige Befehle beschränkt?

- Hardwarekosten: Wie kostspielig ist das Konzept bezüglich der zusätzlich benötigten Hardware?
- Verfügbare Rechenleistung: Wie hoch ist die Verfügbare Rechenleistung im Vergleich zu den anderen Konzepten? Reicht diese aus um das Sprachmodell auszuführen?
- Performanz: Als wie performant wird die Lösung eingeschätzt? Ist mit Latenzen zwischen der Einagbe des Benutzers, der Ausgabe einer Antwort und Ausführung der Aktion zu rechnen?
- Overhead: Wie hoch ist im Allgemeinen der Mehraufwand einzuschätzen?

4.3 Konzept A: Spracherkennung und -verarbeitung mittels Arduino

Ein erstes Konzept sieht vor, dass das Audiosignal direkt von einem der Arduinos in der Getränkemischmaschine aufgenommen wird. Das Audiosignal wird vom Arduino interpretiert und eine passende Antwort wird ausgegeben. Außerdem sendet der Arduino die entsprechenden Signale, um die vom Benutzer gewünschte Aktion von der Getränkemischmaschine ausführen zu lassen. Ein Problem ist hierbei die Interpretation des Audiosignals durch den Arduino, da dessen Leistung nicht für das Ausführen eines Sprachmodells ausreicht. Folglich muss dieser Prozess ausgelagert werden. Das Konzept wird deshalb um ein cloudbasiertes Sprachverarbeitungssystem ergänzt, welches den Sprachbefehl des Benutzers vom Arduino entgegennimmt und einen passenden Befehl und eine passende Antwort zurückgibt (s. Abb. 4.1). Die Kommunikation zwischen Arduino und Cloudsystem kann über das Hypertext Transfer Protocol (HTTP) erfolgen.

4.3. KONZEPT A: SPRACHERKENNUNG UND -VERARBEITUNG MITTELS ARDUINO42

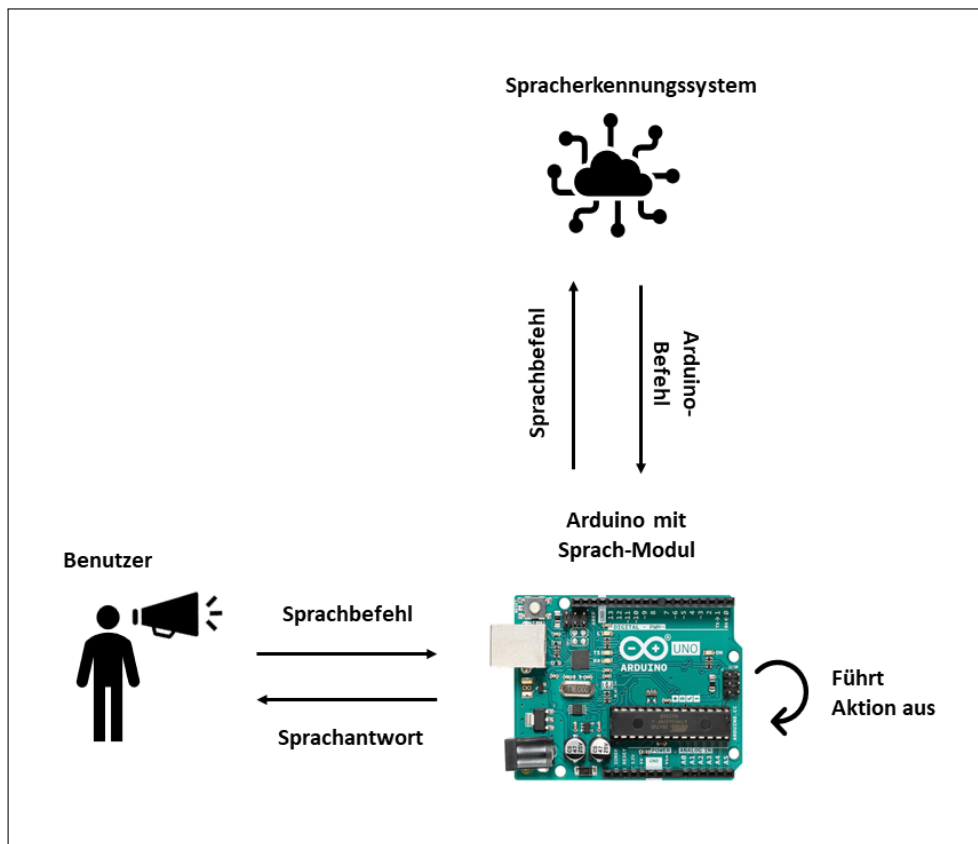


Abbildung 4.1: Spracherkennung und -verarbeitung mittels Arduino

Es muss ein geeignetes Format zum Versenden des Sprachbefehls über HTTP gefunden werden. Eine Möglichkeit besteht darin, das eingehende Audiosignal im Arduino in textform umzuwandeln und diesen String zu versenden. Die auf dem Markt verfügbaren Arduino-Sprach-Module sind jedoch nicht in der Lage beliebige Spracheingaben in Text umzuwandeln, sondern bieten diese Funktionalität nur für vordefinierte Werte an. Dies würde das Ziel dieser Arbeit verfehlen, dem Benutzer eine Konversation mit der Mischmaschine zu ermöglichen. Ein weiteres Problem dieser Lösung besteht darin, dass beispielsweise bei wechselnder Getränkeauswahl die zur Verfügung stehenden Sprachbefehle wie „Ich hätte gerne Getränk xy“ jedes Mal aufs neue manuell angepasst werden müssten. Dies hat zur Folge, dass auch die reine Spracherkennung aus der Mischmaschine ausgelagert werden muss. Ein denkbare Format sind die rohen Audiosignale, die vom Arduino aufgenommen werden.

4.4 Konzept B: Spracherkennung und -verarbeitung mittels mobiler Anwendung

Die Audiosignale über ein Mikrofon in der Mischmaschine aufzunehmen und eine Antwort über einen Lautsprecher auszugeben, so wie es in Konzept A der Fall ist, kann ein Problem darstellen. Zum Einen wird dadurch zusätzliche Hardware benötigt und zum Anderen muss diese korrekt verbaut werden. Das Tonsignal muss vom Mikrofon in einer guten Qualität aufgenommen werden können und die Antwort aus dem Lautsprecher für den Benutzer verständlich sein. Konzept B umgeht dieses Problem durch den Einsatz einer mobilen Anwendung, die durch den Benutzer installiert wird. Über diese Anwendung können anschließend die Aufnahme der Audiosignale, die Spracherkennung und die Kommunikation mit dem Sprachverarbeitungsservice und der Mischmaschine abgewickelt werden, wie in Abbildung 4.2 zu sehen ist.

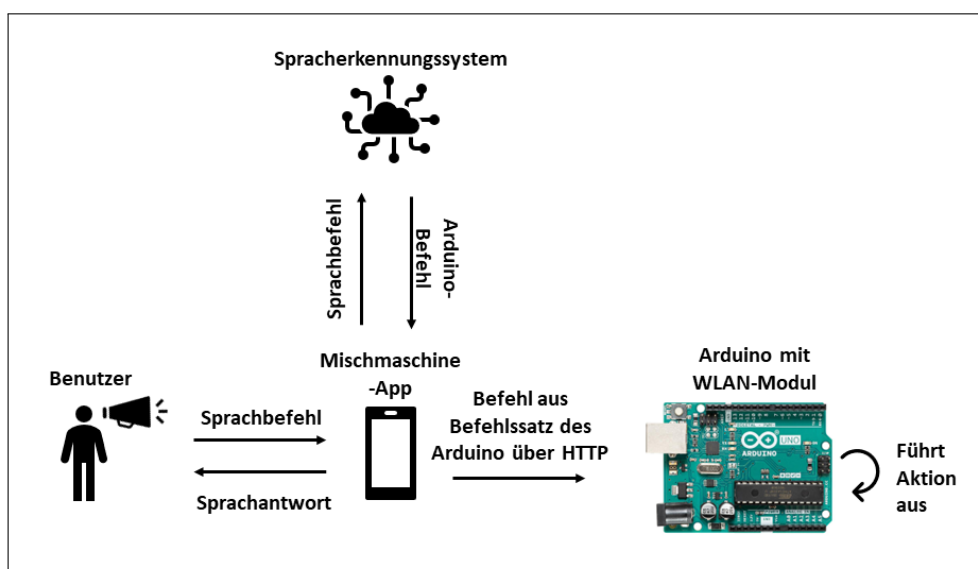


Abbildung 4.2: Spracherkennung und -verarbeitung mittels mobiler Anwendung

Ein Problem dieser Lösung ist der offensichtliche Mehraufwand durch die Entwicklung einer eigenen Anwendung für Mobiltelefone. Auch der Anwender hat zusätzlichen Aufwand durch die Installation. Außerdem ist die Spracheingabe und -ausgabe über das Mobiltelefon nicht intuitiv, da der Anwender eigentlich mit der Maschine kommunizieren sollte. Dieser Effekt kann dadurch abgeschwächt werden, dass wenigstens die Antwort durch einen Lautsprecher in der Mischmaschine an den Benutzer zurückgegeben wird.

4.5 Konzept C: Spracherkennung und -verarbeitung auf Computer-Hardware

Ein weiteres Konzept stützt sich auf die Verwendung eines Computers in der Mischmaschine anstelle eines Mikrocontrollers wie dem Arduino. Motivation ist hierbei der Leistungsgewinn gegenüber eines Mikrocontrollers, um die Spracherkennung und -verarbeitung mittels Sprachmodell zu gewährleisten. Ein Beispiel für einen solchen Miniaturcomputer ist der Raspberry-Pi. Dieser bietet genügend Schnittstellen, wie etwa USB-Hubs, zum Verbinden von Mikrofon als auch Lautsprecher. Nimmt der Computer das Audiosignal auf verarbeitet er dieses und generiert daraus die Antwort, die durch den Lautsprecher ausgegeben wird, zusammen mit der Aktion für die Getränkemischmaschine. Diese muss an den Arduino, welcher die Mischmaschine steuert, übermittelt werden. Um dies zu ermöglichen können der Computer und der Arduino über eine serielle Schnittstelle, wie etwa einem USB-Kabel, miteinander verbunden werden. Abbildung 4.3 stellt den konzeptionellen Aufbau graphisch dar.

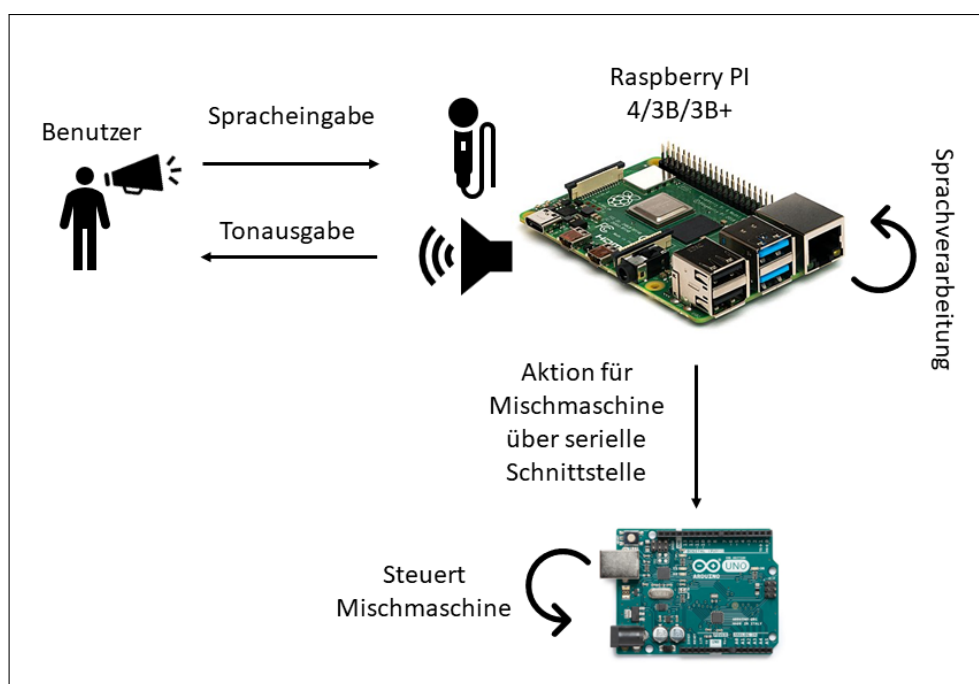


Abbildung 4.3: Spracherkennung und -verarbeitung auf Computer-Hardware

Obwohl ein Computer wie der Raspberry-Pi im Allgemeinen eine höhere Leistung als ein Mikrocontroller hat ist damit nicht sichergestellt, dass diese zur Ausführung des Sprachmodells ausreicht. Beispielsweise ist das vierte Modell der Raspberry-Pi-Serie mit nur maximal acht Gigabyte Arbeitsspeicher erhältlich. Das Sprachmodell könnte allerdings noch weitaus mehr Daten im Arbeitsspeicher benötigen. Des Weiteren ist zu beachten, dass die Miniaturcomputer von Raspberry-Pi im Speziellen zum Zeitpunkt dieser Arbeit kaum zu vertretbaren Preisen

verfügbar sind.

4.6 Finales Hardware-Konzept

Die folgende Tabelle zeigt eine Übersicht bezüglich der Bewertung der einzelnen Konzepte.

Bewertungsmatrix	Konzept A (nur Arduino)	Konzept A	Konzept B	Konzept C
Freiheitsgrade in der Spracheingabe des Benutzers	niedrig	sehr hoch	sehr hoch	hoch
Hardwarekosten	niedrig	hoch	niedrig	sehr hoch
Verfügbare Rechenleistung	niedrig	sehr hoch	sehr hoch	hoch
Performanz	sehr hoch	mittel	mittel	sehr hoch
Overhead	niedrig	hoch	sehr hoch	niedrig

Tabelle 4.1: Bewertung der Konzepte

Ein erster Ansatz wurde im Kapitel 4.3 zu Konzept A erläutert. Hierbei war die Idee, die Spracherkennung und -verarbeitung nur auf dem Arduino auszuführen. Aufgrund der mangelnden Leistung eines Arduinos können mit Hilfe von Sprachmodulen jedoch nur vordefinierte Sätze oder Wörter erkannt werden. Damit ist der Freiheitsgrad in der Spracheingabe des Benutzers äußerst eingeschränkt. Dafür sind die aufzuwendenden Hardwarekosten minimal. Lediglich das Sprachmodul sowie ein Lautsprecher müssten besorgt werden. Die Performanz wird als sehr hoch eingeschätzt, da die Spracherkennung direkt auf dem Arduino erfolgen kann, der auch die Mischmaschine steuert. Niedrig ist hingegen der benötigte Mehraufwand, da kaum zusätzliche Anwendungen, Services oder Hardwarekomponenten benötigt werden.

Das Konzept wurde schließlich durch einen Sprachverarbeitungsservice in der Cloud ergänzt (Konzept A). Der Benutzer hat bei diesem Ansatz große Freiheiten in seinen Formulierungen, da es kein Problem darstellt ein großes Sprachmodell in der Cloud auszuführen. Die Hardwarekosten könnten allerdings hoch sein, je nach dem, ob der Server selbst bereitgestellt oder von einem externen Anbieter bezogen wird. Auch die letzten Endes tatsächlich benötigte Rechenleistung spielt dabei eine Rolle. Die verfügbare Rechenleistung ist theoretisch unbegrenzt, wobei die Performanz des Gesamtsystems nur als mittelmäßig eingestuft werden kann. Nach der Aufnahme und eventuell einer Vorverarbeitung des Audiosignals durch den Arduino müssen HTTP-Nachrichten gesendet und Empfangen werden. Je nach Last auf dem Netzwerk kann es dadurch zu Latenzen oder sogar Verbindungsabbrüchen kommen. Außerdem ist der Overhead durch den Einsatz einer Cloud recht hoch.

Konzept B unterscheidet sich in Sachen Freiheitsgrad, Rechenleistung und Performanz nicht von Konzept A, da auch hier eine Cloud zum Einsatz kommt. Die Hardwarekosten sind jedoch niedriger, da immerhin kein Sprachmodul, Mikrofon und Lautsprecher benötigt werden. Die Aufgaben dieser Komponenten kann das Mobiltelefon des Anwenders übernehmen. Der Overhead ist deutlich größer, da das Konzept die Entwicklung einer eigenen Mobilanwendung voraussetzt.

Der Freiheitsgrad wird bei Konzept C als hoch, jedoch nicht als sehr hoch, bewertet. Grund hierfür ist die, im Vergleich zur Cloud, etwas beschränkte Leistung, welche die Leistung/Größe des Sprachmodells beeinträchtigen könnte. Hardwarekosten können jedoch sehr hoch werden. Bei dem Einsatz einer Cloud kann ein günstiger Anbieter gefunden werden, sodass die Anschaffung eigener Hardware entfällt. Dies ist hier nicht der Fall. Die Performanz des Gesamtsystems kann, wie bei Konzept A (nur mittels Arduino), sehr hoch eingeschätzt werden, da Spracherkennung und -verarbeitung direkt in der Mischmaschine von der Hardware übernommen wird. Der Mehraufwand ist gering, da weder ein Cloudservice noch eine externe Anwendung entwickelt werden müssen.

Letzten Endes wurde das Konzept C als finales Hardware-Konzept gewählt. Wie ein Blick in die Bewertungsmatrix zeigt lässt sich dies mit den sehr guten Werten, die sich aus der Betrachtung der beschriebenen Bewertungskriterien ergaben, begründen.

4.7 Konzept für die Sprachsteuerung

4.7.1 Ansatz für das Dialogsystem

Aus dem Vergleich der wichtigsten Ansätze zur Erstellung von Chatbots, die in Kapitel 2.3.5 besprochen wurden, lassen sich die folgenden Vor- und Nachteile der einzelnen Methoden ableiten:

Ansatz	Vorteile	Nachteile
Musterabgleich	<ul style="list-style-type: none"> • Einfacher Einstieg • Leicht wiederverwendbar • Modularität • Leicht zu kontrollieren/ einzuschränken 	<ul style="list-style-type: none"> • Themenbereich begrenzt • Die Möglichkeiten sind durch die Arbeitsbelastung des Entwicklers begrenzt • Komplexität der Fehlersuche • Strenge und „spröde“ Regeln
Grounding	<ul style="list-style-type: none"> • Gut im Beantworten logischer Fragen • Leicht zu kontrollieren/ einzuschränken 	<ul style="list-style-type: none"> • Künstlicher, mechanischer Ton • Probleme mit Zweideutigkeiten • Probleme mit dem Allgemeinwissen • Begrenzt auf strukturierte Daten • Erfordert die Extraktion von Informationen in großem Umfang • Erfordert menschliche Aufsicht
Suche	<ul style="list-style-type: none"> • Einfachheit • Leicht zu lehren • Simulation von menschlicher Konversation 	<ul style="list-style-type: none"> • Unzureichende Skalierung • Die simulierte Persönlichkeit des Bots ist inkonsistent • Kennt den Kontext nicht • Keine sachlichen Fragen
Generierungsmethoden	<ul style="list-style-type: none"> • Neue, kreative Dialoge • Weniger Arbeit für den Entwickler • Kontextsensitiv 	<ul style="list-style-type: none"> • Schwierig zu lehren • Erfordert mehr Daten (Dialoge) • Schwierig, in die richtige Richtung zu lenken • Erfordert mehr Rechenleistung

Tabelle 4.2: Bewertung der Ansätze für die Erstellung eines Dialogsystems

Bei der Analyse des Problems, ein Dialogsystem für eine Getränkmischmaschine zu entwickeln, kann man zu dem Schluss kommen, dass die beste Option eine Mischung aus dem Ansatz der Informationssuchemethode und dem Musterabgleich ist.

Das erste Argument, das für diesen Ansatz spricht, ist die Möglichkeit, die Sprachsteuerung in deutscher Sprache zu verwenden, was die Anwendung der Generierungsmethoden erschwert, da der Zugang zu einer geeigneten Datenbank in dieser Sprache, die den Anforderungen des Projekts, nämlich eine ausreichende Anzahl von Beleidigungen zu enthalten, nicht möglich ist. Gleichzeitig können bei der Verwendung von Musterabgleich Antwortvorlagen und Muster für

entsprechende Anfragen in deutscher Sprache im Voraus erstellt werden, was die Erstellung und das Training des Dialogsystems erleichtert.

Das zweite Argument ist, dass in diesem Projekt ein Raspberry Pi verwendet wird, was die Verwendung generativer Methoden aufgrund der begrenzten Hardware-Ressourcen einschränken kann. Andererseits kann das Modell für die Informationssuche und den Musterabgleich auf Geräten mit geringem Stromverbrauch implementiert werden, was diesen Ansatz für dieses Projekt vorteilhaft macht.

Generierungsmethoden könnten auch für dieses Problem zu mächtig sein. Beim Mixen von Cocktails können die Antworten einfach und formelhaft sein, wie z. B. „Das hört sich eklig an, bist du sicher, dass du es willst?“ oder „Ich hoffe, ich sehe dich nie wieder“. In diesem Fall kann der Einsatz von Generierungsmethoden wie seq2seq-Modellen überflüssig und ineffizient sein. Für diese Aufgabe ist im Gegensatz zu komplexen natürlichsprachlichen Abfragen keine detaillierte semantische Verarbeitung erforderlich, so dass der Musterabgleich einen einfacheren und effizienteren Ansatz darstellt.

Bei der Verwendung des Musterabgleiches kann man den Ton und den Humor des Geräts leicht steuern und so die gewünschte Atmosphäre erzeugen. Die Antworten des seq2seq-Modells sind wiederum sehr schwer zu steuern. Generierungsmethoden können zu unerwünschtem Maschinenverhalten führen, wenn das Modell auf ungeeigneten Daten trainiert wird oder Fehler in der Betriebslogik enthält. Die Verwendung einer Datenbank, die genügend Beleidigungen enthält (z. B. die 4chan-Datenbank), kann dazu führen, dass die Antworten der Maschine über das Ziel hinausschießen und statt lustig zu sein, den Benutzer beleidigen.

Mit diesem Ansatz wird auch die Effizienz der Maschinensteuerung verbessert. Beim Musterabgleich kann eine Kategorie „Getränkebestellung“ zugewiesen werden, die, wenn sie erkannt wird, die Maschine automatisch zur Bearbeitung der Befehle veranlasst. Im Falle von seq2seq-Modell ist keine Kategorie vorgesehen, so dass eine zusätzliche Prüfung jeder Eingabeanweisung eingeführt werden müsste, um festzustellen, wann der Benutzer die Bestellung aufgegeben hat.

Ein letztes Argument, das für den Ansatz spricht, ist die Möglichkeit, die Maschine bei Bedarf schnell an neue Anfragen und Anforderungen anzupassen, indem neue Muster und Regeln in das System eingeführt werden. Daher ist eine Mischung aus Informationssuchemethode und Musterabgleich für das vorliegende Problem am besten geeignet.

4.7.2 Befehle

Kapitel 5

Implementierung

5.1 Implementierung des Sprachverarbeitungssystems

5.1.1 Word2Vec-Modell

5.1.2 Sequence-to-Sequence-Modell

5.2 Implementierung der Sprachsteuerung

Im Folgenden wird erläutert, wie die Sprachsteuerung für die Getränkemischmaschine implementiert wurde und welche Technologien dafür zum Einsatz kamen. Dabei wird zunächst auf die Spracherkennung d.h., die Umwandlung der Audiosignale (Sprachbefehl des Benutzers) in eine Form, die innerhalb des Quelltextes weiterverarbeitet werden kann, eingegangen (s. Abschnitt 5.2.1). Danach wird die Anbindung des Sprachverarbeitungssystems beschrieben, dessen Implementierung in Abschnitt 5.1 erklärt wird. Abschließend wird die Kommunikation mit der Mischmaschine über den Arduino illustriert (s. Abschnitt 5.2.3).

5.2.1 Spracherkennung

Die Spracherkennung ist der erste Schritt bei der Implementierung einer Sprachsteuerung für die Getränkemischmaschine. Mit Spracherkennung ist die Aufnahme eines Tonsignals über ein Audio-Eingabegerät (Mikrofon) und die Umwandlung der Audiodaten in Text gemeint. Der Quelltext zur Implementierung der Sprachsteuerung erfolgt mit der Programmiersprache *Python*, da hier sehr viele, leicht zu bedienende Bibliotheken zur Spracherkennung, -verarbeitung und KI zur Verfügung stehen.

Für dieses Projekt fiel die Wahl auf das Paket *SpeechRecognition*, das die Verwendung verschiedener Spracherkennungsdienste über eine einheitliche Schnittstelle ermöglicht und zu diesem Zweck auch zur Aufnahme und Verarbeitung der Audiosignale verwendet werden kann [ZHANG (UBER)]

o. D.] Ein großer Vorteil davon ist, dass dadurch ein schneller Wechsel der eingesetzten API erfolgen kann, sollte dies erforderlich sein. Die Verwendung des *SpeechRecognition*-Pakets findet fast ausschließlich über die *Recognizer*-Klasse statt. Um Audiosignale über eine physische Audioquelle (bspw. ein Mikrofon am Computer) aufzunehmen kann die *Microphone*-Klasse verwendet werden, die ebenfalls im Paket enthalten ist. Mit Hilfe eines Objekts vom Typ *Microphone* und der Methode *listen* der *Recognizer*-Klasse können anschließend Audiodaten aufgenommen werden, die in einem Objekt vom Typ *AudioData* gespeichert sind. Die Verwendung von *Recognizer* und *Microphone* sind in Listing 5.1 zu sehen.

```
1 import speech_recognition as sr

3 r = sr.Recognizer()
4 m = sr.Microphone()

6 with m as source:
7     audio = r.listen(source)
```

Listing 5.1: Audioaufnahme mit *SpeechRecognition*

Das *AudioData*-Objekt kann nun verwendet werden um die darin gespeicherten Audiodaten zu erkennen und in Text umzuwandeln. Das *SpeechRecognition*-Paket stellt dafür verschiedene Möglichkeiten zur Verfügung, wie eingangs erwähnt wurde. Diese sollen im Folgenden kurz beschrieben werden:

- **Whisper:** Whisper ist ein neuronales Netz das von der Firma *OpenAI* trainiert und als Open-Source-Projekt zur Verfügung gestellt wird [OPENAI o.D.[a],[b]]. Neben der Fähigkeit Sprache in Text zu konvertieren kann es auch eingesetzt werden um Transkripte zu generieren, die gesprochene Sprache automatisch zu erkennen oder in die englische Sprache zu übersetzen. Das Paket *SpeechRecognition* lässt sowohl die Verwendung der von *OpenAI* zur Verfügung gestellten Online-API zu, als auch die lokale Ausführung des Sprachmodells. Aufgrund der Anforderung nach Offline-Funktionalität entfällt die erste Möglichkeit (s. Kapitel 3).
- **Sphinx:** Das *CMUSphinx* Projekt wird von der Carnegie Mellon University (CMU) unterhalten und stellt eine Reihe von Werkzeugen und Bibliotheken zur Spracherkennung zur Verfügung [SHMYREV o.D.]. Darunter fallen bspw. *PocketSphinx*, *SphinxTrain* und *sphinx4*. Bei *PocketSphinx* handelt es sich um eine C-Bibliothek zur Spracherkennung, die auch innerhalb von Python verwendet werden kann. *SphinxTrain* hingegen stellt Ressourcen zum Trainieren eigener Modelle bereit. *sphinx4* ist das Java-Equivalent zu *PocketSphinx*. Somit ist für dieses Projekt nur *PocketSphinx* interessant. Leider unterstützt enthält das Paket standardmäßig nur ein Modell für die englische Sprache. Anderssprachige Modelle und weitere dazugehörige Ressourcen müssen mühsam aus externen Quellen bezogen und eingebunden werden.
- **Snowboy:** *Snowboy Hotword Detection* ist ein quelloffenes Projekt welches Bibliotheken für verschiedene Programmiersprachen bereitstellt, die dem Erkennen sog. „hot words“

dienen [KIT-AI 2023]. Ein *hot word* sind ein oder mehrere Wörter die der Aktivierung eines Sprachsteuerungssystems dienen. Bekannte Beispiele hierfür sind „Alexa“ oder „OK Google“. Das System „hört“ im Hintergrund und wartet darauf, dass das *hot word* genannt wird. Anschließend nimmt das System weitere Sätze, die vom Benutzer gesprochen werden, wahr und verarbeitet diese. Ein *hot word* ist notwendig damit das System nicht dauerhaft versucht Sätze zu verarbeiten, die nicht für das System gedacht sind und damit wiederholt Fehler bekommt oder unerwünschtes Verhalten an den Tag legt. Da die Ansprüche des vorliegenden Projekts an die Detektion von *hot words* nicht besonders hoch sind und um die von Snowboy benötigte Rechenzeit einzusparen wird *Snowboy* in diesem Projekt keine Relevanz haben. Auf die konkrete Umsetzung des *hot word* Mechanismus wird später noch eingegangen. Ein weiterer Grund dafür, dass Snowboy nicht zum Einsatz kommen soll ist, dass das Projekt bereits im Jahr 2020 offiziell eingestellt wurde und nur noch von der Community weiter betrieben wird.

- Vosk: Das *Vosk Speech Recognition Toolkit* ist neben *Whisper* und *Sphinx* die dritte und letzte verfügbare offline API für die *SpeechRecognition* eine Schnittstelle bereitstellt [ALPHA CEPHEI o. D., 2023]. Es wird von der *Alpha Cephei Incorporation* vertrieben. Ähnlich wie bei *Whisper* sind vortrainierte Modelle in vielen verschiedenen Sprachen bereits vorhanden, müssen allerdings manuell heruntergeladen und dem *SpeechRecognition* Paket zur Verfügung gestellt werden.
- Google: *SpeechRecognition* bietet zwei Möglichkeiten Google-Dienste zur Konvertierung von Sprache zu Text zu verwenden. Die erste Möglichkeit besteht in der Verwendung der *Google Speech API* in der *Google Cloud* [GOOGLE o. D.] Die zweite Möglichkeit besteht in der Verwendung der *Google Speech API* ohne die *Google Cloud*. Die Verwendung der Google-Schnittstelle entfällt jedoch für dieses Projekt, da eine aktive Internetverbindung Voraussetzung dafür ist. Außerdem sind die Google-Dienste weder kostenfrei noch open-source.
- Microsoft: Auch hier gibt es zwei verschiedene Möglichkeiten Microsoft-Dienste zur Umwandlung von Sprache zu Text über das *SpeechRecognition*-Paket zu nutzen. Eine davon ist die Verwendung der *Speech to text* Funktion in der *Microsoft Azure Cloud* [MICROSOFT o. D.] Hierfür sind ein aktiver *Azure*-Account sowie ein gültiger API-Schlüssel notwendig. Die zweite Möglichkeit besteht in der *Microsoft Bing Voice Recognition*. Diese ist jedoch veraltet und wird nicht mehr unterstützt. Auch diese Schnittstelle entfällt jedoch wegen der notwendigen Internetverbindung und Bezahlung.
- IBM: Auch die Firma IBM bietet ihre eigene Schnittstelle zur Spracherkennung mit Hilfe der von IBM entwickelten KI *Watson* an [IBM 2023]. Auch hier gilt: zur Verwendung der API ist sowohl eine aktive Internetverbindung als auch Bezahlung vorgesehen, weshalb diese Schnittstelle für das Projekt nicht in Frage kommt.

- Weitere Schnittstellen sind *Wit.ai*, *Houndify* und *Tensorflow* [SOUNDHOUND o.D.; TENSORFLOW o.D.[a]; WIT.AI o.D.] Da diese Schnittstellen allderdings ebenso eine aktive Internetverbindung voraussetzen und damit ausscheiden soll nicht weiter auf sie eingegangen werden.

Die einzelnen APIs lassen sich jeweils über einen Methodenaufruf der Form *Recognizer.recognize_x* verwenden, wobei das „x“ für den Namen der jeweiligen API steht. Für dieses Projekt viel die Wahl auf die Verwendung der *Whisper*-Bibliothek von *OpenAI*. Dies ist damit zu begründen, dass sie von den offline verwendbaren APIs die mit Abstand am einfachsten zu verwendende ist und gleichzeitig sehr gute Ergebnisse beim Testen damit erzielt wurden. Beispielsweise werden die Modelle für viele verschiedene Sprachen bereits von OpenAI zur Verfügung gestellt und für jedes Modell stehen weitere Ausführungen zur Verfügung die nach den eigenen Ansprüchen und vorhandenen Ressourcen ausgewählt werden können. Die verschiedenen Ausführungen sind nach der „Größe“ des Modells unterteilt in „tiny“, „base“, „small“, „medium“ und „large“ [OPENAI 2023].

```

1 import speech_recognition as sr

3 r = sr.Recognizer()
4 m = sr.Microphone()

6 with m as source:
7     print("Start listening...")
8     audio = r.listen(source)
9     text = ""
10    try:
11        recognized_text = r.recognize_whisper(audio, language="german", model="tiny")
12        text = recognized_text
13        print("Recognized text: " + text)
14    except sr.UnknownValueError:
15        print("Whisper could not understand audio.")
16    except sr.RequestError as e:
17        print("Could not request results from Whisper; {0}".format(e))

```

Listing 5.2: Sprache zu Text mit *OpenAI Whisper*

Listing 5.2 zeigt, wie mit Hilfe der *Whisper*-API die aufgenommenen Audiodaten zu Text verarbeitet und ausgegeben werden können. Die Sprache und die Modellgröße werden bei dem Methodenaufruf *Recognizer.recognize_whisper* angegeben. Sollte das entsprechende Modell noch nicht lokal vorliegen wird dieses bei der ersten Ausführung automatisch installiert, was einen hohen Grad an Benutzerfreundlichkeit seitens der API bedeutet. Natürlich ist es unzureichend einmal die Sprache des Benutzers aufzunehmen und in Text zu konvertieren. Insgesamt muss das Programm während des Betriebs der Getränkemischmaschine in einer Dauerschleife ausgeführt werden und der Benutzer muss die Möglichkeit haben mehrere Befehle nacheinander über die Sprachsteuerung auszuführen. Listing 5.3 zeigt wie dies konkret erreicht wurde.

```

1 # imports and error handling

3 r = sr.Recognizer()
4 m = sr.Microphone()
5 with m as source:
6     r.adjust_for_ambient_noise(source)

8 while True:

```



```

9      asound.snd_lib_error_set_handler(c_error_handler)
10     print("Start listening ...")
11     with m as source:
12         audio = r.listen(source)
13     asound.snd_lib_error_set_handler(None)
14     text = ""
15     try:
16         recognized_text = r.recognize_whisper(audio, language="german", model="tiny")
17         text = recognized_text
18         print("Recognized text: " + text)
19     except sr.UnknownValueError:
20         print("Whisper could not understand audio.")
21     except sr.RequestError as e:
22         print("Could not request results from Whisper; {0}".format(e))
23
24     if text.lower().strip().startswith("mischmaschine"):
25         print("Recognized hot word. Will output answer and command.")

```

Listing 5.3: Hot-Word-Erkennung

Der Code der für die Konvertierung in Sprache zu Text zuständig wird nun innerhalb einer While-True-Schleife ausgeführt. Um zu verhindern, dass die Mischmaschine Befehle verarbeitet, die keine sind, soll das *hot word* „Mischmaschine“ zur Aktivierung genutzt werden. Deshalb wird in Zeile 24 nun geprüft, ob der erkannte Text mit diesem Wort beginnt und nur dann werden weitere Schritte zur Befehlsverarbeitung eingeleitet (mehr dazu in Abschnitt 5.2.3). Da davon auszugehen ist, dass die Mischmaschine an Orten eingesetzt wird, die vielen Umgebungsgeräuschen ausgesetzt sind wird in Zeile sechs außerdem die Methode `speech_recognition.Recognizer().adjust_for_ambient_noise(source)` aufgerufen, die dafür sorgt, dass die Empfindlichkeit des Mikrofons automatisch an die Lautstärke der Umgebungsgeräusche angepasst wird.

Herausforderungen bei der Implementierung der Spracherkennung

Die Verwendung von *Whisper* setzt die Installation von *PyTorch* voraus - einem quelloffenen ML-Framework [PYTORCH o.D.] Die Installation von *PyTorch* kann auf zwei unterschiedlichen Wegen erfolgen:

1. Unter der Benutzung vorkompilierter Distributionen mit Hilfe eines Paketverwaltungssystems wie etwa *pip* [PPA o.D.]
2. Durch das eigenständige Bauen des Quelltextes auf dem lokalen Rechner.

Da es sich bei *PyTorch* um ein überaus großes und komplexes Projekt handelt ist es empfehlenswert die erste Methode anzuwenden, da der Build-Vorgang bei Methode zwei sehr viele Ressourcen benötigt.

Nach der Verwendung von *pip* zur Installation von *PyTorch* auf dem für das Projekt verwendeten *Raspberry Pi 4B* kam es jedoch zu einem unerwarteten Fehler. Das Importieren der *PyTorch*-Bibliothek ließ das Programm sofort beenden und die Nachricht „Illegal Instruction“ auf dem Bildschirm ausgeben. Die Fehlermeldung lässt darauf schließen, dass die Bibliothek einen

für den Prozessortyp unbekannten Befehl ausüben will. Drei Lösungsansätze wurden daraufhin konzipiert:

1. Das Programm sollte mit einer anderen Python-Version ausgeführt werden.
2. Das vorkompilierte Paket, das automatisch von *pip* ausgewählt wurde scheint nicht auf die vorhandene Prozessorarchitektur zu passen und führt unbekannte Maschinenbefehle aus. Deshalb sollte die *PyTorch*-Bibliothek selbstständig, lokal kompiliert werden.
3. Ein anderes vorkompiliertes Paket herunterladen und installieren, das auf die vorhandene Prozessorarchitektur angepasst und optimiert ist.
4. Sollten alle zuvor genannten Versuche scheitern muss *Whisper* durch eine andere Bibliothek ersetzt werden, die keine Abhängigkeit zu *PyTorch* aufweist.

Lösungsansatz eins erwies sich als nicht effektiv. Das Programm verursachte immernoch den selben Fehler. Lösungsansatz zwei scheiterte, da der Build-Prozess mehr als den zur Verfügung stehenden RAM des Raspberry Pi aufbrauchte. Sogar mit vier Gigabyte (GB) zusätzlich konfiguriertem Swap-Speicher wurde der Build-Prozess nach einiger Zeit vom Betriebssystem beendet aufgrund des hohen RAM-Verbrauchs. Lösungsansatz drei funktionierte schließlich. Mit Hilfe des vorkompilierten Pakets aus [Q-ENGINEERING o. D.] konnte *PyTorch* installiert und erfolgreich importiert und verwendet werden.

Ein weiteres Problem trat bei der Verwendung des Mikrofons über die *SpeechRecognition*-Bibliothek auf. Bei den Aufrufen *speech_recognition.Microphone()* und *speech_recognition.Recognizer().listen(source)* kommt es auf dem Raspberry Pi zu unerwarteten Fehlermeldungen, die in Abbildung 5.1 zu sehen sind.

```
ALSA lib conf.c:4745:(_snd_config_evaluate) function snd_func_refer returned error: No such file or directory
ALSA lib conf.c:5233:(snd_config_expand) Evaluate error: No such file or directory
ALSA lib pcm.c:2660:(snd_pcm_open_noupdate) Unknown PCM front
ALSA lib confmisc.c:1281:(snd_func_refer) Unable to find definition 'cards.bcm2835_headpho.pcm.surround51.0:CARD=0'
ALSA lib conf.c:4745:(_snd_config_evaluate) function snd_func_refer returned error: No such file or directory
ALSA lib conf.c:5233:(snd_config_expand) Evaluate error: No such file or directory
ALSA lib pcm.c:2660:(snd_pcm_open_noupdate) Unknown PCM surround21
ALSA lib confmisc.c:1281:(snd_func_refer) Unable to find definition 'cards.bcm2835_headpho.pcm.surround51.0:CARD=0'
ALSA lib conf.c:4745:(_snd_config_evaluate) function snd_func_refer returned error: No such file or directory
ALSA lib conf.c:5233:(snd_config_expand) Evaluate error: No such file or directory
ALSA lib pcm.c:2660:(snd_pcm_open_noupdate) Unknown PCM surround21
ALSA lib confmisc.c:1281:(snd_func_refer) Unable to find definition 'cards.bcm2835_headpho.pcm.surround40.0:CARD=0'
ALSA lib conf.c:4745:(_snd_config_evaluate) function snd_func_refer returned error: No such file or directory
ALSA lib conf.c:5233:(snd_config_expand) Evaluate error: No such file or directory
ALSA lib pcm.c:2660:(snd_pcm_open_noupdate) Unknown PCM surround40
ALSA lib confmisc.c:1281:(snd_func_refer) Unable to find definition 'cards.bcm2835_headpho.pcm.surround51.0:CARD=0'
ALSA lib conf.c:4745:(_snd_config_evaluate) function snd_func_refer returned error: No such file or directory
```

Abbildung 5.1: Fehlermeldungen beim Zugriff auf das Mikrofon

Zwar kann das Programm trotz der Meldungen fehlerfrei ausgeführt werden, jedoch stören die Meldungen bspw. beim Debuggen und irritieren evtl. Programmierer, die zu einem späteren Zeitpunkt einen Blick auf das Projekt werfen. Um die Fehlermeldungen zu umgehen wurde deshalb im finalen Quelltext die Fehlerbehandlungsfunktion der *libasound*-Bibliothek ausgetauscht wie in Listing 5.4 zu sehen ist.

```

1 from ctypes import *
2 import speech_recognition as sr

4 ERROR_HANDLER_FUNC = CFUNCTYPE(None, c_char_p, c_int, c_char_p, c_int, c_char_p)
5 def py_error_handler(filename, line, function, err, fmt):
6     pass
7 c_error_handler = ERROR_HANDLER_FUNC(py_error_handler)

9 asound = cdll.LoadLibrary("libasound.so")
10 asound.snd_lib_error_set_handler(c_error_handler)

12 r = sr.Recognizer()
13 m = sr.Microphone()

15 with m as source:
16     asound.snd_lib_error_set_handler(c_error_handler)
17     print("Start listening...")
18     audio = r.listen(source)
19     asound.snd_lib_error_set_handler(None)
20     text = ""
21     try:
22         recognized_text = r.recognize_whisper(audio, language="german", model="tiny")
23         text = recognized_text
24         print("Recognized text: " + text)
25     except sr.UnknownValueError:
26         print("Whisper could not understand audio.")
27     except sr.RequestError as e:
28         print("Could not request results from Whisper; {0}".format(e))

```

Listing 5.4: Fehlerbehandlung

Für einen Vorher-/Nachher-Vergleich sei auf Listing 5.2 verwiesen. Zeilen vier bis zehn dienen der Definition einer C-Funktion zur Fehlerbehandlung, die bei ihrer Ausführung die auftretende Meldung ignoriert. In Zeile 16 wird die Standard Fehlerfunktion durch diese ersetzt und in Zeile 19 wieder zurückgeändert.

5.2.2 Anbindung des Sprachmodells an die Mischmaschine

5.2.3 Befehlsverarbeitung in der Mischmaschine

Nachdem die Sprache des Benutzers erkannt, in Text konvertiert und durch das Sprachmodell interpretiert wurde muss der zurückerhaltene Befehl vom Raspberry Pi and den Arduino in der Getränkemischmaschine gesendet werden. Zur Umsetzung der Kommunikation zwischen Raspberry Pi und Arduino kommt die *PySerial*-Bibliothek zum Einsatz [PYSERIAL o. D.] Listing 5.5 zeigt, wie die Bibliothek zum Versenden von Nachrichten an ein Gerät über die serielle Schnittstelle erfolgen kann.

```

1 import serial
2 import time

4 ser = serial.Serial("/dev/ttyACM0", 9600, timeout=5)
5 ser.reset_input_buffer()
6 time.sleep(5)

```

```

7 count = 0
8 while True:
9     ser.write("Message {0}\n".format(count).encode("utf-8"))
10    line = ser.readline().decode("utf-8").rstrip()
11    print(line)
12    count += 1
13    time.sleep(1)

```

Listing 5.5: Serielle Kommunikation in Python

Das Programm öffnet in Zeile vier die serielle Schnittstelle zu dem Gerät mit dem Namen „ttyACM0“, mit einer Baud-Rate von 9600 und einem Timeout von fünf Sekunden. Der Gerätenamen ist je nach Gerät und Betriebssystem unterschiedlich und muss zuvor mit den verfügbaren Betriebssystemmitteln bestimmt werden. Die Baud-Rate gibt an, mit welcher Geschwindigkeit Symbole über die Schnittstelle gesendet werden. Die Baud-Rate muss sowohl beim Sender als auch beim Empfänger gleich sein, da es sonst zu Kommunikationsproblemen kommt. Der Timeout steuert wie lange bei einer Leseoperation gewartet werden soll, bis alle zu lesenden Bytes im Puffer vorhanden sind.

In der darauffolgenden Zeile wird der Eingabepuffer geleert, um zu verhindern, dass unvollständige Daten an den Empfänger gesendet werden, die beim letzten Programmdurchlauf im Puffer verblieben sein könnten. Da der Verbindungsaufbau mit dem Gerät einige Zeit dauern kann wird die Programmausführung in Zeile sechs zunächst um fünf Sekunden verzögert. Das nachfolgende Programm sendet in einer While-True-Schleife mehrere Nachrichten über die serielle Schnittstelle, die über einen Zähler voneinander zu unterscheiden sind. Anschließend wird in Zeile zehn eine Nachricht zurück erwartet die anschließend auf der Kommandozeile ausgegeben wird.

In Listing 5.6 ist der Quelltext auf Empfängerseite zu sehen, der dafür zuständig ist die in Listing 5.5 generierten Nachrichten zu empfangen und Daten zurückzusenden.

```

1 void setup() {
2     Serial.begin(9600);
3 }
4
5 void loop() {
6     if (Serial.available() > 0) {
7         String data = Serial.readStringUntil('\n');
8         Serial.println(data);
9         delay(1000);
10    }

```

Listing 5.6: Serielle Kommunikation in Python

Das Programm ist in der Arduino-Programmiersprache verfasst und besteht, wie vorgegeben, aus einer *setup* und einer *loop* Methode. Die *setup* Methode wird einmalig beim Start des Programms, wenn der Arduino mit Strom versorgt wird, ausgeführt wohingegen die *loop* Methode dauerhaft während der Arduino mit Strom versorgt ist wiederholt ausgeführt wird. In der *setup*-Methode wird die serielle Schnittstelle mit einer Baud-Rate von 9600 geöffnet. Innerhalb der *loop*-Methode wird geprüft, ob Bytes in den Eingabepuffer geschrieben wurden. Ist dies der Fall wird der

If-Block betreten und die in den Puffer geschriebenen Bytes werden bis zu einem vorgegebenen Terminalsymbol (in diesem Fall ein Zeilenumbruch) eingelesen und in ein String-Objekt konvertiert. Mit dem Aufruf *Serial.println(data)* wird der empfangene String über die serielle Schnittstelle zurück an den ursprünglichen Sender übertragen. Dabei handelt es sich um den String, der in dem Programm aus Listing 5.5 in Zeile zehn erwartet wird.

Die eben an einem einfachen Beispiel beschriebene Vorgehensweise/Implementierung kann eins zu eins auf die Problemstellung des vorliegenden Projekts übertragen werden. Statt der generischen Nachricht wird das Kommando, das von dem Sprachmodell zurückgegeben wird übertragen. Auf Seiten des Arduinos muss das Kommando interpretiert und die entsprechenden Aktionen innerhalb der Mischmaschine ausgelöst werden. Listing 5.7 zeigt wie auf seiten des Raspberry Pi die eben beschriebenen Techniken in das Hauptprogramm aus Listing 5.4 integriert werden können.

```

1 from ctypes import *
2 import speech_recognition as sr
3 import serial
4 import time

6 ERROR_HANDLER_FUNC = CFUNCTYPE(None, c_char_p, c_int, c_char_p, c_int, c_char_p)
7 def py_error_handler(filename, line, function, err, fmt):
8     pass
9 c_error_handler = ERROR_HANDLER_FUNC(py_error_handler)

11 asound = cdll.LoadLibrary("libasound.so")
12 asound.snd_lib_error_set_handler(c_error_handler)

14 def language_model_stub(s):
15     return s, "50,50,0,0\n"

17 r = sr.Recognizer()
18 m = sr.Microphone()

20 ser = serial.Serial("/dev/ttyACM0", 9600, timeout=5)
21 ser.reset_input_buffer()
22 time.sleep(5)

24 with m as source:
25     asound.snd_lib_error_set_handler(c_error_handler)
26     print("Start listening...")
27     audio = r.listen(source)
28     asound.snd_lib_error_set_handler(None)
29     text = ""
30     try:
31         recognized_text = r.recognize_whisper(audio, language="german", model="tiny")
32         text = recognized_text
33         print("Recognized text: " + text)
34     except sr.UnknownValueError:
35         print("Whisper could not understand audio.")
36     except sr.RequestError as e:
37         print("Could not request results from Whisper; {0}".format(e))

39     answer, cmd = language_model_stub(text)
40     ser.write(cmd.encode("utf-8"))
41     serial_answer = ser.readline().decode("utf-8").rstrip()
42     print(serial_answer)

```

Listing 5.7: Serielle Kommunikation für die Mischmaschine

In den Zeilen 14 und 15 wird eine Methode definiert, die das Sprachmodell, welches die Sprache des Benutzers interpretiert und einen Maschinenbefehl und eine Antwort zurückgibt, nachahmen

soll. Das Kommando besteht aus vier natürlichen Zahlen, die die Prozentwerte der jeweiligen Behälter in der Mischmaschine darstellen. Der String „50,50,0,0“ ist also so zu interpretieren, dass das Mischungsverhältnis zu 50 Prozent aus Behälter eins und zu 50 Prozent aus Behälter zwei bestehen soll. In Zeile 40 wird dieses Kommando an den Arduino, welcher die Mischmaschine steuert, übertragen. Listing 5.8 zeigt einen Ausschnitt des abgeänderten Code auf Seiten des Arduinos.

```

1  int getValueAt(String s, int vpos) {
2      int pos = 0;
3      int valueCount = 0;
4      for (int i = 0; i < s.length(); i++) {
5          if (s.charAt(i) == ',' || i == s.length() - 1) {
6              String value = "";
7              int bound = (i == s.length() - 1 ? i + 1 : i);
8              for (int j = pos; j < bound; j++) {
9                  value.concat(s.charAt(j));
10             }
11             if (valueCount == vpos) {
12                 return value.toInt();
13             }
14             pos = i + 1;
15             valueCount++;
16         }
17     }
18     return -1;
19 }

21 void setup() {
22     Serial.begin(9600);
23 }

25 void loop() {
26     if (Serial.available() > 0) {
27         String data = Serial.readStringUntil('\n'); // e.g. "50,50,0,0"
28         Serial.println("Received " + data);
29         analogWrite(pumpe1, (getValueAt(data, 0) * 0.5 + 50) * 2.55);
30         analogWrite(pumpe2, (getValueAt(data, 1) * 0.5 + 50) * 2.55);
31         analogWrite(pumpe3, (getValueAt(data, 2) * 0.5 + 50) * 2.55);
32         analogWrite(pumpe4, (getValueAt(data, 3) * 0.5 + 50) * 2.55);
33         delay(3000); // let the "button" be "pressed" for three seconds
34     } else {
35         analogWrite(pumpe1, 0);
36         analogWrite(pumpe2, 0);
37         analogWrite(pumpe3, 0);
38         analogWrite(pumpe4, 0);
39         analogWrite(pumpe5, 0);
40     }
41 }

```

Listing 5.8: Interpretation der Kommandos

Mit der Methode *getValueAt(String, int)*: *int* wird die Prozentangabe an der gewünschten Stelle im Kommandostring ausgelesen. In der Hauptprogrammschleife werden die Pumpen eins bis vier mit ihren entsprechenden Werten angesprochen. Das Delay in der letzten Zeile des If-Blocks sorgt dafür, dass die Pumpen nicht sofort wieder in der nächsten Schleifeniteration ausgeschaltet werden.

Herausforderungen bei der Implementierung der Befehlsverarbeitung

Eine Herausforderung bestand in der Steuerung der seriellen Kommunikation von Raspberry Pi und Touch-Display mit dem Arduino.

Kapitel 6

Fazit und Ausblick

TODO

Literatur

AIML Foundation [o. D.] URL: <http://www.aiml.foundation/doc.html> [besucht am 12.03.2023] [siehe S. 28].

ALAMMAR, Jay [o. D.] *Sequence-to-Sequence Models for Chatbots* [siehe S. 33].

ALPHA CEPHEI [o. D.] *VOSK Offline Speech Recognition API*. en. URL: <https://alphacephei.com/vosk/> [besucht am 03.04.2023] [siehe S. 51].

— [Apr. 2023]. *Vosk Speech Recognition Toolkit*. original-date: 2019-09-03T17:48:42Z. URL: <https://github.com/alphacep/vosk-api> [besucht am 03.04.2023] [siehe S. 51].

AMAZON [o. D.] *Amazon Lex: How It Works - Amazon Lex*. URL: <https://docs.aws.amazon.com/lex/latest/dg/how-it-works.html> [besucht am 12.03.2023] [siehe S. 31].

CHAWLA, Sumit [o. D.] *Building Chatbots with Google Dialogflow: Create chatbots with Dialogflow's natural language processing and machine learning capabilities* [siehe S. 31].

DEERWESTER, Scott u. a. [1990]. »Indexing by latent semantic analysis«. en. In: *Journal of the American Society for Information Science* 41.6. __eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/%28199009%2941%3A6%3C391%3A%3AAID-ASI1%3E3.0.CO%3B2-9>, S. 391–407. ISSN: 1097-4571. DOI: 10.1002/(SICI)1097-4571(199009)41:6<391::AID-ASI1>3.0.CO;2-9. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/%28SICI%291097-4571%28199009%2941%3A6%3C391%3A%3AAID-ASI1%3E3.0.CO%3B2-9> [besucht am 06.04.2023] [siehe S. 23].

DIANA, Perez-Marin und Pascual-Nieto ISMAEL [Juni 2011]. *Conversational Agents and Natural Language Interaction: Techniques and Effective Practices: Techniques and Effective Practices*. en. Google-Books-ID: 2nUcqtBCOBcC. IGI Global. ISBN: 978-1-60960-618-3 [siehe S. 28, 30].

ECKART, Carl und Gale YOUNG [Sep. 1936]. »The approximation of one matrix by another of lower rank«. en. In: *Psychometrika* 1.3, S. 211–218. ISSN: 1860-0980. DOI: 10.1007/BF02288367. URL: <https://doi.org/10.1007/BF02288367> [besucht am 06.04.2023] [siehe S. 23, 24].

- GOLDBERG, Yoav [Apr. 2017]. *Neural Network Methods in Natural Language Processing*. Englisch. Hrsg. von Graeme HIRST. San Rafael, Calif.: Morgan & Claypool Publishers. ISBN: 978-1-62705-298-6 [siehe S. 11].
- GOOGLE [o. D.] *Speech-to-Text: Automatic Speech Recognition*. en. URL: <https://cloud.google.com/speech-to-text> [besucht am 26.03.2023] [siehe S. 37, 51].
- GÖTZ, Felix, Moritz HÖCKELE und Florian LOBERT [2022]. *Mischmaschine*. Projektarbeit Software des Studiengangs Mechatronik an der DHBW Karlsruhe [siehe S. 1, 3].
- HERCZEG, Michael [März 2018]. »9. Zeitverhalten interaktiver Systeme«. de. In: *9. Zeitverhalten interaktiver Systeme*. De Gruyter Oldenbourg, S. 173–182. ISBN: 978-3-11-044686-9. DOI: 10.1515/9783110446869-187. URL: <https://www.degruyter.com/document/doi/10.1515/9783110446869-187/html> [besucht am 26.03.2023] [siehe S. 36].
- IBM [März 2023]. *IBM Watson Speech to Text - Overview*. en-us. URL: <https://www.ibm.com/cloud/watson-speech-to-text> [besucht am 02.04.2023] [siehe S. 51].
- ISELL, Charles und Paul VIOLA [1998]. »Restructuring Sparse High Dimensional Data for Effective Retrieval«. In: *Advances in Neural Information Processing Systems*. Bd. 11. MIT Press. URL: https://proceedings.neurips.cc/paper_files/paper/1998/hash/87ec2f451208df97228105657edb717f-Abstract.html [besucht am 06.04.2023] [siehe S. 23].
- JANSEN, Jonas [Feb. 2017]. »Digitalisierung: 8,4 Milliarden vernetzte Geräte im Internet der Dinge«. de. In: *FAZ.NET*. ISSN: 0174-4909. URL: <https://www.faz.net/aktuell/wirtschaft/netzwirtschaft/digitalisierung-8-4-milliarden-vernetzte-geraete-im-internet-der-dinge-14865654.html> [besucht am 15.03.2023] [siehe S. 1].
- JURAFSKY, Dan u. a. [2009]. *Speech and language processing: an introduction to natural language processing, computational linguistics, and speech recognition / Daniel Jurafsky (Stanford University), James H. Martin (University of Colorado at Boulder)*. eng. ISBN: 978-0-13-504196-3 [siehe S. 5, 22].
- KITT-AI [Apr. 2023]. *Snowboy Hotword Detection*. original-date: 2016-05-10T00:54:30Z. URL: <https://github.com/Kitt-AI/snowboy> [besucht am 03.04.2023] [siehe S. 51].
- KOECH, Kiprono Elijah [Juli 2022]. *Cross-Entropy Loss Function*. en. URL: <https://towardsdatascience.com/cross-entropy-loss-function-f38c4ec8643e> [besucht am 06.04.2023] [siehe S. 35].
- LANE, Hobson, Cole HOWARD und Hannes Max HAPKE [2019]. *Natural language processing in action: understanding, analyzing, and generating text with Python / Hobson Lane, Cole Howard, Hannes Max Hapke*. eng. ISBN: 978-1-61729-463-1 [siehe S. 6, 27].

- LI, Jiwei u. a. [Sep. 2017]. *Adversarial Learning for Neural Dialogue Generation*. arXiv:1701.06547 [cs]. DOI: 10.48550/arXiv.1701.06547. URL: <http://arxiv.org/abs/1701.06547> [besucht am 12. 03. 2023] [siehe S. 32].
- LTD, RASPBERRY PI [o. D.] *Raspberry Pi 4 Model B specifications*. en-GB. URL: <https://www.raspberrypi.com/products/raspberry-pi-4-model-b/specifications/> [besucht am 24. 03. 2023] [siehe S. 38].
- MARWEDEL, Peter [2021]. *Eingebettete Systeme: Grundlagen Eingebetteter Systeme in Cyber-Physikalischen Systemen*. de. Wiesbaden: Springer Fachmedien. ISBN: 978-3-658-33436-9 978-3-658-33437-6. DOI: 10.1007/978-3-658-33437-6. URL: <https://link.springer.com/10.1007/978-3-658-33437-6> [besucht am 15. 03. 2023] [siehe S. 1].
- MICROSOFT [o. D.] *Speech to Text – Audio to Text Translation / Microsoft Azure*. en-US. URL: <https://azure.microsoft.com/en-us/products/cognitive-services/speech-to-text> [besucht am 02. 04. 2023] [siehe S. 51].
- NUPUR SHARMA, Anupriya [o. D.] *Chatbot Development using Machine Learning Techniques* [siehe S. 32].
- OPENAI [o. D.[a]]. *About*. en-US. URL: <https://openai.com/about> [besucht am 01. 04. 2023] [siehe S. 50].
- [o. D.[b]]. *Introducing Whisper*. en-US. URL: <https://openai.com/research/whisper> [besucht am 01. 04. 2023] [siehe S. 50].
- [Apr. 2023]. *Whisper*. original-date: 2022-09-16T20:02:54Z. URL: <https://github.com/openai/whisper> [besucht am 02. 04. 2023] [siehe S. 52].
- PPA, Python Packaging Authority [o. D.] *pip: The PyPA recommended tool for installing Python packages*. URL: <https://pip.pypa.io/> [besucht am 03. 04. 2023] [siehe S. 53].
- PYSERIAL [o. D.] *Welcome to pySerial's documentation — pySerial 3.0 documentation*. URL: <https://pythonhosted.org/pyserial/> [besucht am 04. 04. 2023] [siehe S. 55].
- PYTORCH [o. D.] *PyTorch*. en. URL: <https://www.pytorch.org> [besucht am 03. 04. 2023] [siehe S. 53].
- Q-ENGINEERING [o. D.] *Install PyTorch on Raspberry Pi 4 - Q-engineering*. en-GB. URL: <https://qengineering.eu/install-pytorch-on-raspberry-pi-4.html> [besucht am 04. 04. 2023] [siehe S. 54].
- RYGL, Jan u. a. [Juni 2017]. *Semantic Vector Encoding and Similarity Search Using Fulltext Search Engines*. arXiv:1706.00957 [cs]. URL: <http://arxiv.org/abs/1706.00957> [besucht am 06. 04. 2023] [siehe S. 21].
- SCIKIT-LEARN [o. D.] *scikit-learn documentation — DevDocs*. en. URL: https://devdocs.io/scikit_learn/ [besucht am 12. 03. 2023] [siehe S. 30].

- SHMYREV, Nickolay [o. D.] *About CMUSphinx*. URL: <http://cmusphinx.github.io/wiki/about/> [besucht am 02.04.2023] [siehe S. 50].
- SOUNDHOUND [o. D.] *Home*. en-US. URL: <https://www.soundhound.com/> [besucht am 02.04.2023] [siehe S. 52].
- TENSORFLOW [o. D.[a]]. *TensorFlow*. en. URL: <https://www.tensorflow.org/> [besucht am 02.04.2023] [siehe S. 52].
- [o. D.[b]]. *TensorFlow - Sequence-to-Sequence Models*. URL: <https://chromium.googlesource.com/external/github.com/tensorflow/tensorflow/+r0.7/tensorflow/g3doc/tutorials/seq2seq/index.md> [besucht am 12.03.2023] [siehe S. 33].
- WIT.AI [o. D.] *Wit.ai*. URL: <https://wit.ai/> [besucht am 02.04.2023] [siehe S. 52].
- WONG, S. K. M., Wojciech ZIARKO und Patrick C. N. WONG [Juni 1985]. »Generalized vector spaces model in information retrieval«. In: *Proceedings of the 8th annual international ACM SIGIR conference on Research and development in information retrieval*. SIGIR '85. New York, NY, USA: Association for Computing Machinery, S. 18–25. ISBN: 978-0-89791-159-7. DOI: 10.1145/253495.253506. URL: <https://dl.acm.org/doi/10.1145/253495.253506> [besucht am 06.04.2023] [siehe S. 23].
- WOUDENBERG, Aswin van [2014]. »A Chatbot Dialogue Manager-Chatbots and Dialogue Systems: A Hybrid Approach«. Magisterarb. Open Universiteit Nederland [siehe S. 27].
- ZHANG (UBERI), Anthony [o. D.] *SpeechRecognition: Library for performing speech recognition, with support for several engines and APIs, online and offline*. URL: https://github.com/Uberi/speech_recognition#readme [besucht am 30.03.2023] [siehe S. 49].