

Information Retrieval

Seminararbeit

des Studienganges Angewandte Informatik
an der Dualen Hochschule Baden-Württemberg Mannheim

von

Jesse-Jermaine Richter, Jonas Seng

29.04.2019

Matrikelnummer, Kurs:	8787549/1980179, TINF16AIBI
Ausbildungsfirma:	DZ BANK AG, Frankfurt
Betreuer der Ausarbeitung:	Herr Prof. Dr. Karl Stroetmann

Erklärung

Wir versichern hiermit, dass wir unsere Seminararbeit mit dem Thema: „Information Retrieval“ selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt haben.

Wir versichern zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

Ort, Datum

Unterschrift

Ort, Datum

Unterschrift

Die Seminararbeit mit dem Thema „Information Retrieval“ hat als Ziel, dem Leser das Thema Information Retrieval näher zu bringen und dieses anhand einer Beispielimplementierung in Python praktisch zu veranschaulichen. Hierbei wird in den ersten Kapiteln das theoretische Wissen an den Leser vermittelt, damit dieser die Beispielimplementierung in Kapitel 5 nachvollziehen kann. Dies beinhaltet die Themen Tokenization, invertierter Index und das TF-IDF-Maß. Die Beispielimplementierung wird mit Python umgesetzt und in Kapitel 5 genau beschrieben. Am Ende in Kapitel 6 werden Probleme und Verbesserung der Beispielimplementierung angesprochen.

The goal of this paper with the topic „Information Retrieval“ is to bring the topic Information Retrieval closer to the reader and to illustrate this practically by means of an example implementation in Python. In the first chapters the theoretical knowledge is explained to the reader so that he can understand the example implementation in chapter 5. This includes the topics tokenization, inverted index and the TF-IDF measure. The example implementation is implemented with Python and described in detail in chapter 5. At the end of chapter 6 problems and improvements of the example implementation are discussed.

Inhaltsverzeichnis

1	Einleitung	1
1.1	Was ist Information-Retrieval?	1
1.2	Ziel der Arbeit	2
1.3	Stand der Forschung	2
2	Information Retrieval - Theoretische Grundlagen	4
2.1	Problemstellung	4
2.2	Strategiefindung	5
2.3	Tokenization	6
3	Invertierter Index	9
3.1	Grundlegender Aufbau	9
3.2	Umsetzung eines invertierten Index	9
4	TF-IDF	13
4.1	Term Frequency	13
4.2	Inverse Document Frequency	14
4.3	Das TF-IDF-Maß	15
5	Implementierung - Lokale Suchmaschine	17
5.1	Ziel der Beispiel-Implementierung	17
5.2	Genutzte Bibliotheken	17
5.3	Die Document-Klasse	20
5.4	Der Index	22
5.5	Ausführung der Suchmaschine	30
6	Verbesserungen	32
6.1	Probleme der aktuellen Beispiel-Implementierung	32
6.2	Aufbaudauer des Index	34

Abbildungsverzeichnis

3.1	Beispiel für einen invertierten Index [1]	10
3.2	Beispiel eines Tries [11]	11
5.1	Imports	19
5.2	Dokumentenklasse	20
5.3	TF-IDF-Methode	21
5.4	Indexklasse	22
5.5	BuildIndex-Methode	24
5.6	GetStartDirectories-Methode	25
5.7	AddToIndex-Methode	26
5.8	PreprocessText-Methode	29
5.9	Retrieve-Methode	31
5.10	Execute	31

Abkürzungstabelle

Abkürzung:	Bedeutung:
ASCII	American Standard Code for Information Interchange
DF	Document Frequency
IDF	Inverse Document Frequency
IR	Information Retrieval
PDF	Portable Document Format
TF	Term Frequency
TF-IDF	TF-IDF-Maß

1 Einleitung

1.1 Was ist Information-Retrieval?

Information-Retrieval (IR) beschreibt das Bereitstellen spezieller Informationen aus einer großen und unsortierten Datenmengen. Dieses Themengebiet fällt unter Informatik, Informationswissenschaften sowie Computerlinguistik und ist ein wesentlicher Bestandteil von Suchmaschinen wie zum Beispiel Google [13].

Das Thema besitzt bereits seit einigen Jahren eine hohe, aber dennoch steigende Relevanz. Die Gründe der hohen Relevanz von IR liegen vor allem beim Einsatz von Suchmaschinen. Diese sind in Zeiten des Internets die wohl wichtigste Form der Informationsbeschaffung. Aufgrund der immer schneller steigenden Informationsmengen wird das Thema künftig weiter an Relevanz gewinnen. Unternehmen und Privatanwendern wird eine immer wachsende Menge von Informationen zugänglich, die organisiert werden muss, damit relevante bzw. spezifisch gesuchte Informationen jederzeit und ohne Verzögerung gefunden werden.

Um das Ziel der Bereitstellung von Informationen gewährleisten zu können, wird erst eine Durchsuchung und Gewichtung sämtlicher Informationen bzw. Dokumente, die später gefunden werden sollen, durchgeführt. Das zentrale Objekt der Informationsrückgewinnung stellt der invertierte Index dar [4, S. 6], dessen Aufbau und Funktionsweise in den nächsten Kapiteln ausführlich erläutert wird. Im Verlauf dieser Arbeit wird das TF-IDF-Maß, welches zur Beurteilung der Relevanz eines Dokumentes genutzt wird, im Fokus stehen.

Die theoretischen Hintergründe des invertierten Index und des TF-IDF-Maßes werden durch eine Beispiel-Implementierung einer lokalen Suchmaschine in der Programmiersprache Python veranschaulicht.

1.2 Ziel der Arbeit

Ziel der Arbeit ist es, ein grundlegendes Verständnis des Themenkomplexes Information-Retrieval zu vermitteln. Das umfasst auch die theoretischen Hintergründe, die für die später vorgestellten Beispiel-Implementierungen notwendig sind.

Die Beispiel-Implementierung soll hauptsächlich die folgenden Themengebiete umfassen:

- Tokenization von Texten (Kapitel 2.3 und 5.4.2)
- Aufbau eines invertierten Indexes (Kapitel 3 und 5.4)
- Approximierende Beurteilung der Relevanz eines gefundenen Dokuments mittels TF-IDF (Kapitel 4 und 5.3.1)
- Mögliche Verbesserungen der Implementierung (Kapitel 6)

Die in dieser Arbeit vorgestellte Implementierung hat nicht den Anspruch auf eine hohe Performance, vielmehr dient diese dem Zwecke der praxisnahen Veranschaulichung der Funktionsweise von IR-Systemen.

1.3 Stand der Forschung

Dieser Abschnitt umreißt kurz den aktuellen Stand der Forschung. Dazu werden zwei Modelle von Information Retrieval knapp beschrieben, die für die Entwicklung einer lokalen Suchmaschine, von Bedeutung sind. Es wird jedoch nur auf ein Modell im Verlauf dieser Arbeit eingegangen.

1.3.1 Vector Space Model

Das Vector Space Model, zu deutsch Vektorraummodell, repräsentiert Dokumente und Anfragen als hochdimensionale, metrische Vektoren [6]. Der Anfrage-Vektor wird beim Retrieval-Prozess mit den Dokumenten-Vektoren verglichen. Dabei werden jedoch nur Do-

kumente betrachtet, welche mit der Anfrage in Verbindung stehen [12]. Welche Dokumente betrachtet werden, wird mithilfe des invertierten Index ermittelt.

Es gibt verschiedene Maße, mit denen die Vektoren miteinander verglichen werden. Der einfachste Ansatz besteht darin, den Abstand zwischen Anfrage-Vektor und Dokumenten-Vektor zu berechnen. Weiter verbreitet ist das Cosinus-Maß (oder Cosinus-Ähnlichkeit), welches den Winkel zwischen Anfrage-Vektor und Dokumenten-Vektor angibt. Je kleiner der Winkel zwischen den Vektoren, desto höher ist die Relevanz des Dokuments [2].

Dieses Modell wird im Rahmen dieser Arbeit genauer beleuchtet und als Grundlage für die Beispiel-Implementierung genutzt.

2 Information Retrieval - Theoretische Grundlagen

2.1 Problemstellung

Im Rahmen dieser Arbeit soll beschrieben werden, wie mithilfe von „Information Retrieval“ Informationen aus einer großen, unsortierten Datenmenge - nach Relevanz sortiert - bereitgestellt werden kann. Dabei bekommt das System eine vom Nutzer gestellte Abfrage, auch Query genannt, und versucht auf deren Basis, Daten, die meist als Dokumente vorliegen, zurückzuliefern. Im Gegensatz zu Abfragen im Datenbankumfeld beinhaltet die Query jedoch unzureichende Informationen, um ein spezielles Element eindeutig identifizieren zu können [4, S. 1 f.]. Dies soll ein IR-System auch nicht leisten. Vielmehr sollen Ergebnisse zurückgeliefert werden, die mit hoher Wahrscheinlichkeit Relevanz bzgl. der gestellten Query besitzen. Der Nutzer selektiert dann die für diesen nötigen Dokumente [4, S. 1 f.]. Mathematisch lässt sich dies folgendermaßen formulieren: Aus einer Dokumentenmenge D soll mithilfe einer Funktion eine Teilmenge D_1 von D ermittelt werden, die relevant für eine Abfrage q ist.

Um diese Funktion sinnvoll definieren zu können, muss jedoch zuvor die Menge aller Queries, sowie die Menge aller Tokens definiert werden. Zuvor sollen jedoch die Vokabeln „Token“, „Typ“ und „Term“ geklärt werden:

Definition 2.1 (Token) *Unter einem Token wird eine zusammenhängende Zeichenkette verstanden, die innerhalb eines Dokuments vorkommt [4, S. 22].*

Tokens werden häufiger auch als „Wörter“ bezeichnet, dies wird auch im Rahmen dieser Arbeit so gehandhabt.

Definition 2.2 (Typ) *Ein Typ bezeichnet eine Klasse von Tokens, die dieselben Zeichen in derselben Reihenfolge enthalten [4, S. 22].*

Definition 2.3 (Term) *Ein Term ist ein Typ, welcher im Dictionary eines IR-Systems vorkommt [4, S. 22].*

Ein Beispiel, wie Tokens, Typen und Terme aussehen können, folgt in 2.3.2.

Definition 2.4 (Menge aller Terme) *Sei $d \in D$ ein Dokument. Die Menge T_d ist nun die Menge aller Wörter, die in dem Dokument d enthalten sind: $T_d = \{t_1, \dots, t_n\}$.*

Die Menge T ist die Menge aller Terme, die in den Dokumenten aus D vorkommen, also: $T = T_{d_1} \cup \dots \cup T_{d_i}$ mit $i \in N$ und $d_i \in D$.

Mithilfe der Definition 2.4 kann nun die Menge aller möglichen Queries definiert werden:

Definition 2.5 (Menge aller möglichen Queries) $Q \subseteq 2^T$

Definition 2.6 (Retrievalfunktion) *Eine Funktion $f: Q \rightarrow 2^D$ heißt Retrievalfunktion, wobei Q die Menge aller Queries ist.*

Nachdem die Problemstellung formuliert ist, muss eine Strategie entwickelt werden, wie die Retrievalfunktion nach Definition 2.6 dargestellt bzw. umgesetzt wird.

2.2 Strategiefindung

Dieser Abschnitt bietet eine Übersicht, wie das (im Folgenden vorgestellte) IR-System arbeiten soll.

Als Vorarbeit werden alle Dokumente, die im Index aufgenommen werden, in eine Codierung wie ASCII oder Unicode umgewandelt. Dazu wird ein Tool genutzt, das hier nicht weiter von Relevanz sein wird. Es sollen mindestens all diejenigen Dokumente in den Index aufgenommen werden, die im PDF-Format vorliegen.

Der erste Schritt, der das IR-System an sich leisten muss, ist das Erstellen von Tokens (siehe Definition 2.1). Dazu wird jedes Dokument in Tokens aufgespalten. Ein Token ist in den meisten Fällen ein Wort. Satzzeichen wie Punkte, Kommata usw. sollen nicht als

Tokens behandelt und dementsprechend ignoriert werden, dies gilt ebenso für Leerzeichen und Tabs.

Für jedes Token wird es später im Index einen Eintrag geben, der eine Liste mit weiteren Informationen hält. Diese Liste muss mindestens die Dokument-IDs speichern, in welchen das Token steht. In diesen Listen werden häufig noch weitere Informationen hinterlegt, beispielsweise die Häufigkeit eines Tokens.

Der zweite Schritt besteht darin, einen Algorithmus zu entwerfen, der eine Query entgegennimmt und auf Basis der Query und des Index eine Liste von relevanten Dokumenten ausgibt. Dieser Algorithmus wird das in der Einleitung kurz vorgestellte Vektorraummodell verwenden. Weiter wird dieser für die Ermittlung der Relevanz die sogenannte TF-IDF-Gewichtung nutzen. Diese wird später noch ausführlich vorgestellt.

2.3 Tokenization

2.3.1 Vorarbeiten

Bevor aus Dokumenten Tokens erzeugt werden können, müssen einige Fragen beantwortet werden. Eine Frage ist, welche Dokumente betrachtet werden sollen und wie man ein Dokument definiert. Ein Beispiel soll das Problem veranschaulichen:

Angenommen das IR-System soll dazu dienen Dokumente auf der Festplatte eines Computers zu finden. In diesem Szenario kann jede in einem Ordner gelistete Datei als Dokument angesehen werden. Dies wäre der einfachste Fall. Jedoch ist dies meist nicht erwünscht. So sollen beispielsweise bestimmte Dateitypen von der Suche ausgeschlossen werden. In UNIX existiert ein Dateityp, welcher mehrere Mails pro Datei speichert. Hier wird jede Mail als einzelnes Dokument angesehen. Daraus folgt, dass die Maildatei in mehrere Dokumente aufgespalten werden muss [4]. Umgekehrt gibt es Szenarien, in denen mehrere Dokumente zu einem Dokument zusammengefasst werden müssen, um bei der Suche nutzbare Ergebnisse zu erzielen [4]. Auch hier dienen E-Mails wieder als Beispiel: In vielen Mails werden Anhänge versandt. Häufig sind dies PDF-Dateien und Dateien von Text-Bearbeitungsprogrammen. Die Anhänge sollen zu dem E-Mail-Dokument gezählt werden, an dem sie angehängt wurden.

Ein weiteres Problem, das gelöst werden muss, um die Dokumente verarbeiten zu können, ist die Codierung der Inhalte der Dokumente. Hierbei müssen Dokumente, die meist in vielen unterschiedlichen Codierungen vorliegen, zu einer definierten Codierung überführt werden [4].

Diese Probleme der „Vorarbeit“ werden in der später gezeigten Beispielimplementierung nicht behandelt, dies wird von anderen Tools übernommen.

2.3.2 Tokenerzeugung

Sobald geklärt ist, in welcher einheitlichen Codierung die Dokumente vorliegen und was als Dokument, im Englischen auch „document unit“ genannt, verstanden wird, kann ein Dokument in Tokens aufgeteilt werden.

Eine wichtige Frage, die im Rahmen der Tokenerzeugung geklärt werden muss, ist, welche Zeichenketten als Token behandelt werden. Kommata, Punkte und sonstige Satzzeichen haben keine sinnvolle Bedeutung im Zusammenhang mit Information Retrieval. Diese Zeichen können somit aus Tokens entfernt bzw. während der Tokenerzeugung überlesen werden [4]. Nach der Definition 2.1, erzeugt der Text

Beispielsatz, der ein Komma und ein Punkt hat.

die Tokens:

[Beispielsatz, der, ein, Komma, und, ein, Punkt, hat]

Die dazugehörige Menge von Typen sieht nach Definition 2.2 so aus:

$$Types = \{Beispielsatz, der, ein, Komma, und, Punkt, hat\}$$

Die erzeugte Menge von Termen könnte nach Definition 2.3 - je nach Verarbeitung der Typen - wie folgt aussehen:

$$Terms = \{beispielsatz, der, ein, komma, und, punkt, hat\}$$

Einige Information Retrieval-Systeme nutzen darüber hinaus sogenannte „stop words“. Das sind Wörter, die in sehr vielen Dokumenten in großer Anzahl vorkommen und damit wenig Bedeutung für die Suche besitzen [4]. Beispiele für solche Wörter sind „ist“, „sein“ und „und“. Jedoch funktioniert diese Technik später beim Suchen schlechter als zunächst

angenommen. Das Wort „sein“ kann beispielsweise als Verb oder als Pronomen in einem Dokument vorkommen. Als Pronomen kann dieses Wort durchaus wichtig sein für eine Suche (beispielsweise innerhalb eines Buchtitels), wird jedoch als stop word aussortiert. Daher werden in neuen IR-Systemen entweder gar keine stop words oder nur eine geringe Anzahl stop words genutzt [4, S. 27].

Eine weitere Möglichkeit solche Wörter zu filtern, ist „Stemming“. Diese Methode führt Wörter auf ihren Wortstamm zurück [4] [3]. Dadurch wird die Anzahl der Terme, die im Index gespeichert werden, stark gesenkt. Allerdings bringt diese Methode eine Unschärfe mit sich. Damit ist gemeint, dass zwei nicht verwandte Wörter auf denselben Wortstamm zurückgeführt werden, wodurch bei der späteren Suche nach einem der beiden Ursprungswörter auch Ergebnisse zurückgeliefert werden, die irrelevant für die Query sind [4].

Weiter existieren sprachspezifische Probleme. Beispielsweise wird im Englischen häufig mit Kurzformen von Wörtern gearbeitet, so wird „are not“ zu „aren’t“. Es muss geklärt werden wie mit solchen Formen umgegangen werden soll. Ein Ansatz ist Query Preprocessing. Hierbei werden Wörter dieser Art, die in der Query und in den Dokumenten stehen, in eine einheitliche Form gebracht [4]. Darüber hinaus gilt es zu beachten, dass Eigennamen wie „Hewlett-Packard“ nicht oder nur nach bestimmten Regeln prozessiert werden dürfen. Welche Wörter als Eigenname behandelt werden und welche nicht, kann mittels Machine Learning-Verfahren oder auf Basis eines großen Vokabulars gelöst werden.

3 Invertierter Index

Dieser Abschnitt wird die Idee des invertierten Index vorstellen sowie die Erzeugung des invertierten Index erläutern. Darüber hinaus wird aufgezeigt, wie die Verarbeitung einer Query mittels des invertierten Index funktioniert. Ab jetzt werden die Ausdrücke „invertierter Index“ und „Index“ synonym verwendet.

3.1 Grundlegender Aufbau

Der invertierte Index kann als ein Dictionary betrachtet werden, welches für jeden Term, der in der Dokumentenmenge vorkommt, einen Eintrag hält [4, S. 6 f.] [3]. Dieser Eintrag wiederum ist eine Liste. Diese hält mindestens eine eindeutige Dokumenten-ID, meist jedoch darüber hinaus weitere Informationen. Beispiele für solche weiteren Informationen sind Häufigkeit, in der ein Term in einem Dokument vorkommt, die Position und die umliegenden Wörter in der Nähe zum Term t . Die Listen, die zu jedem Term angelegt werden, heißen Posting-Listen [4, S. 6].

In Abbildung 3.1 ist ein invertierter Index dargestellt.

3.2 Umsetzung eines invertierten Index

Nachdem das Prinzip des invertierten Index klar ist, steht die Frage im Raum wie dieser umgesetzt werden kann. Es liegt nahe als Datenstruktur einen Trie einzusetzen. Ein Trie ist ein spezieller Suchbaum, welcher besonders gut zum Suchen von Zeichenketten geeignet ist [4, S. 50 f.].

Nummer	Term	Dokumente
1	als	5
2	auch	3
3	das	1, 2, 6
4	der	1, 2, 5
5	dritter	3
6	ein	3
7	ende	6
8	er	4
9	erste	1, 5
10	folgt	3
11	hier	2, 6
12	ist	1, 2, 6
13	länger	4
14	satz	1, 3
15	und	2, 5, 6
16	wobei	4
17	zweite	2, 5

Abbildung 3.1: Beispiel für einen invertierten Index [1]

Alternativ kann über den Einsatz einer Hashmap nachgedacht werden, jedoch führt dies zu folgendem Problem: Angenommen die Terme einer vom Nutzer eingegebenen Query sind nicht in der Hashmap vorhanden, dann wird die Hash-Funktion keinen passenden Eintrag finden. Da in einer Hashmap Wörter, die ähnlich zueinander sind, nicht unbedingt benachbart gespeichert werden und keinerlei Information darüber bekannt ist, wo zur Query ähnliche Wörter gespeichert sind, kann im Falle, dass ein oder mehrere Wörter der Query nicht vorhanden sind, nicht mit geringem Aufwand nach ähnlichen Wörtern gesucht werden. Bei Tries besteht dieses Problem nicht [4, S. 50].

Da die Beispiel-Implementierung, die später vorgestellt wird, jedoch möglichst einfach sein soll, wird zunächst eine Hashmap als Datenstruktur verwendet.

3.2.1 Tries

Ein Trie wird auf der Basis einer Menge von Zeichenketten aufgebaut. Jede Zeichenkette, die gefunden werden muss, ist innerhalb des Tries repräsentiert.

Erreicht wird dies dadurch, dass ein Knoten jeweils ein Zeichen repräsentiert und eine Liste mit Verweisen auf die nächsten möglichen Knoten hält, basierend auf einem weiteren Zeichen. Die Abbildung 3.2 zeigt einen Trie.

Im Folgenden soll eine formale Definition eines Tries gegeben werden:

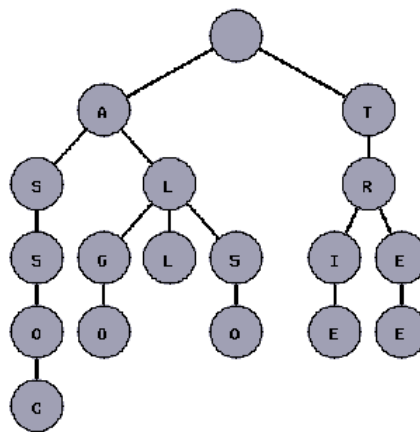


Abbildung 3.2: Beispiel eines Tries [11]

Definition 3.1 (Trie) [7] Sei Σ eine endliche Menge von Zeichen (Alphabet) und Σ^* die Menge aller Wörter, die über Σ gebildet werden können und sei $S \subseteq \Sigma^*$. Sei V die Menge aller Knoten und E die Menge aller Kanten. Jeder Knoten $v \in V$ besitzt einen Wert w , eine Liste $C = [c_1, \dots, c_n]$, wobei $c_1, \dots, c_n \in \Sigma$, und eine Liste $K = [k_1, \dots, k_n]$ von Knoten. Jede Kante $e \in E$ ist eine Verbindung von zwei Knoten $v_1, v_2 \in V$. Dann ist $T = (V, E)$ ist ein Trie, wenn gilt:

- $\forall e \in E : e$ ist mit Zeichen aus Σ beschriftet.
- $\forall v \in V : \text{alle ausgehenden Kanten von } v \text{ sind unterschiedlich beschriftet mit einem Zeichen } z \in \Sigma$
- $\forall S_i \in S : \exists v \in V : S_i \text{ ist ein Präfix der Konkatenation der Beschriftungen des Pfades vom Wurzelknoten bis } v.$
- $\forall b \in V : \exists S_i \in S : \text{Die Konkatenation der Beschriftungen von der Wurzel bis } b \text{ ergibt } S_i, \text{ sofern } b \text{ ein Blatt des Tries ist.}$

Das Suchen nach gespeicherten Wörtern gestaltet sich nun verhältnismäßig einfach: Um das Wort „Tree“ im, in der Abbildung gezeigten, Trie zu finden, wird wie folgt vorgegangen:

Stellt der User die Anfrage „Tree“, wird diese Query nun Zeichen für Zeichen durchgegangen. Beginnend bei „T“ wird im Wurzelknoten geprüft, ob es einen Verweis auf einen

Knoten gibt, der ein „T“ repräsentiert. Existiert ein solcher Knoten, wird in diesem geprüft, ob es einen Knoten gibt, der das nächste Zeichen in der Query (das „r“) repräsentiert. Ist dies der Fall, wird das Verfahren solange wiederholt, bis die Query komplett eingelesen ist oder in der Query ein Zeichen steht, das durch keinen Knoten im Trie repräsentiert ist [11] [5].

Ähnlich leicht funktioniert das Einfügen neuer Wörter in den Trie. Dazu wird - wie beim Suchen - das Wort, das eingefügt werden soll, so weit wie möglich nach dem oben beschriebenen Muster eingelesen und es wird zu den entsprechenden Knoten gesprungen [5]. Wird nun ein Zeichen eingelesen, das nicht durch einen Knoten repräsentiert ist, wird ein neuer Knoten erzeugt, welcher dieses Zeichen repräsentiert. Alle nun noch einzulesenden Zeichen erhalten einen neuen Knoten, da der neu erzeugte Knoten natürlich nicht auf bereits vorhandene Knoten zeigen kann. Innerhalb dieses Knotens wird eine Liste angelegt, die Verweise auf weitere Knoten hält [11]. In dieser Liste wird ein Verweis auf den Knoten angelegt, der das nächste Zeichen des einzufügenden Wortes repräsentiert. Dieses Verfahren setzt sich solange fort, bis das neue Wort vollständig eingelesen ist.

Das Löschen soll in diesem Rahmen nicht aufgezeigt werden, da dies weitaus komplexer sein kann als das Finden oder Einfügen von Einträgen.

Ein Knoten, zu dem man mit dem Wort S_i gelangt, muss außerdem eine Liste mit Verweisen auf alle Dokumente, in denen das Wort S_i vorkommt, speichern. Diese Liste kann nicht leer sein, da das entsprechende Wort S_i nicht im Trie vorkommen würde, sofern es in keinem Dokument steht.

4 TF-IDF

Bei großen Sammlungen von Dokumenten reicht es nicht aus, wenn das IR-System die Dokumente zurückgibt, welche die Suchanfrage erfüllen. Die Menge der Dokumente, die durch das IR zurückgegebene werden, ist meist zu groß, so dass der Nutzer nicht in der Lage ist alle Dokumente zu sichten und die für ihn relevanten Dokument auszuwählen.

Die Frage, die sich hier stellt, ist, wie der Nutzer die Dokumente bekommt, die er mit hoher Wahrscheinlichkeit benötigt. Zur Bestimmung der Wichtigkeit eines Dokuments bzgl. einer Query, wird das sogenannte Scoring eingesetzt. Scoring, oder auch Bewertung, wird genutzt um zu bestimmen, welche Dokumente für die Suchanfrage am relevantesten sind. Es kann auch von einer Gewichtung der Dokumente gesprochen werden.

Das Gewichtungsmodell, welches in dieser Arbeit thematisiert und genutzt wird, ist das TF-IDF-Maß. Das TF-IDF-Maß bewertet ein Dokument anhand der Terme, die in der Query stehen, ohne zu beachten, wie diese miteinander verknüpft sind. TF steht hierbei für Term Frequency und IDF für Inverse Document Frequency. Beide Methoden werden einzeln in den folgenden Abschnitten vorgestellt und am Ende zum TF-IDF-Maß verknüpft.

4.1 Term Frequency

Um Dokumente aus einer Kollektion nach Termen einer Suchquery zu gewichten, ist das Zählen jedes einzelnen Terms t einer Query in jedem Dokument d die zunächst trivialste Lösung. Denn ein Dokument, in dem ein Term häufiger vorkommt als in einem anderes Dokument, weist eine höhere Übereinstimmung mit diesem Term auf.

Bei dieser Vorgehensweise wird von Term Frequency (siehe Definition 4.1) gesprochen, auf deutsch auch als Suchworddichte bekannt.

Definition 4.1 (Term Frequency) [4, S. 117] Die Term Frequency $tf_{d,t}$ gibt die Anzahl des Terms t in dem Dokument d wieder.

Bei dieser Vorgehensweise wird jedoch weder die Anordnung der Terme im Dokument, noch die Anordnung der Terme in der Query berücksichtigt. Die Query „Karl ist engagierter als Eckhard“ liefert die gleiche Gewichtung für die Dokumente einer Kollektion zurück wie „Eckhard ist engagierter als Karl“. Dieses Modell, das die Anordnung der Terme innerhalb des Dokumentes ignoriert, wird auch als „Bag of Words Model“ bezeichnet. Dieses interpretiert ein Dokument als eine Menge von Tupeln $\langle w_i, k \rangle$, wobei w_i das i -te Wort der Menge W ist, welche alle Wörter des Dokuments d beinhaltet. k die Häufigkeit, in der das Wort im Dokument vorkommt. Ein Dokument ist somit ein „Beutel“ mit vielen losen Wörtern. Daraus folgt, dass zwei Dokumente, welche die selben Wörter und die selbe Anzahl dieser haben, als identisch angesehen werden, obwohl sie das nicht unbedingt sind. Des Weiteren werden alle Terme als gleich wichtig gewertet, was nicht immer zielführend ist, wie zum Beispiel im Abschnitt 2.3.2 mit den „stop words“ beschrieben wurde.

4.2 Inverse Document Frequency

Da die Terme in einer Suchquery nicht alle von gleicher Bedeutung für die Suche sind, muss eine Lösung gefunden werden, die die Terme untereinander gewichtet. Wenn in einer Dokumenten-Kollektion im Automobilbereich nach „elektrische Autos sind die Zukunft“ gesucht wird, taucht das Wort Auto in allen Dokumenten häufig auf. Deshalb sollte der Term „Auto“ weniger stark ins Gewicht fallen als das Wort „elektrische“. Hier kommt die Inverse Document Frequency (siehe Definition 4.3) zum tragen. Jedoch muss vorher geklärt werden worum es sich bei der Document Frequency (siehe Definition 4.2) handelt.

Definition 4.2 (Document Frequency) [4, S. 118] Die Document Frequency df_t gibt die Anzahl der Dokumente d in der genutzten Kollektion an, welche den Term t besitzen.

Durch die Document Frequency kann also festgestellt werden, wie viele Dokumente einen bestimmten Term beinhalten. Die Häufigkeit, in der ein Term in den Dokumenten vorkommt, bleibt dabei unbeachtet, dafür ist die Term Frequency zuständig. Um aus der Document Frequency die einzelnen Terme sinnvoll zu gewichten, wird die Inverse Document Frequency gebildet.

Definition 4.3 (Inverse Document Frequency) [4, S. 118] *Die Inverse Document Frequency ist definiert als:*

$$\text{idf}_t = \log\left(\frac{N_D}{\text{df}_t + 1}\right)$$

df_t ist hierbei die Document Frequency und N_D die Anzahl der Dokumente in der Kollektion.

Das Dividieren der Gesamtzahl aller Dokument in der Kollektion (N_D) durch die Document Frequency hat zur Folge, dass das Gewicht eines seltenen Terms höher ist und das eines häufigen Terms niedrig. Der Logarithmus sorgt dafür, dass das resultierende Gewicht nicht zu groß wird, wenn die Kollektion sehr groß und der Term selten ist.

4.3 Das TF-IDF-Maß

In diesen Abschnitt wird die Zusammensetzung der Term Frequency und der Inverse Document Frequency zum TF-IDF-Maß beschrieben.

Definition 4.4 (TF-IDF) [4, S. 119] *Das TF-IDF-Maß ist wie folgt definiert:*

$$\text{tf-idf}_{t,d} = \text{tf}_{t,d} * \text{idf}_t$$

Hierbei wird lediglich das Produkt der Inverse Document Frequency mit der Term Frequency gebildet.

Nun kann eine score-Funktion definiert werden, die einem Dokument ein Gewicht zuweist. Die folgende Funktion $\text{score}(q, d)$ zeigt die Berechnung des Gewichts eines Dokumentes, wobei q hier die Suchquery ist und d das zu gewichtende Dokument: [4, S. 119]

$$\text{score}(q, d) = \sum_{t \in q} \text{tf-idf}_{t,d}$$

Wie in der Funktion beschrieben, wird für jeden Term in der Suchquery das TF-IDF-Maß berechnet. Die einzelnen Gewichte werden dann aufsummiert. Am Ende dieser Prozedur erhält man dann eine Gewichtung für ein Dokument anhand der Suchquery und ist in der Lage seine Relevanz mit der Relevanz der anderen gefundenen Dokumente zu vergleichen.

5 Implementierung - Lokale Suchmaschine

5.1 Ziel der Beispiel-Implementierung

Im Folgenden wird eine Beispiel-Implementierung der zuvor theoretisch diskutierten Inhalte vorgestellt. Dabei wird eine lokale Suchmaschine entwickelt, welche in der Lage ist, PDF-Dateien auf einem lokalen Computer-System zu parsen, in einen invertierten Index aufzunehmen, sowie Suchanfragen eines Benutzers sinnvoll zu beantworten. Zur Relevanz-Bestimmung der Dokumente wird das TF-IDF-Maß, welches bereits vorgestellt wurde, genutzt. Das Speichern des Index wird mit der von Python mitgelieferte Datenstruktur *Dictionary*, welche im Grunde eine Hashmap ist, umgesetzt. Weiter werden Bibliotheken eingesetzt, welche einige Vorarbeit leisten und damit den Code der Beispiel-Implementierung auf das Wesentliche beschränken. Die Anwendung soll die grundlegende Arbeitsweise eines Information Retrieval-Systems darlegen.

5.2 Genutzte Bibliotheken

Vor der eigentlichen Implementierung der lokalen Suchmaschine werden einige Module eingebunden (siehe Abbildung 5.1), welche die Implementierung unterstützen. Im Folgenden werden die Module aufgelistet und wichtige Module im jeweiligen Abschnitt genauer erläutert:

- Apache Tika: Erkennt und extrahiert Metadaten und Texte aus über tausend Dateitypen [8]
- Math: Standard Python-Modul für mathematische Funktionen
- OS: Stellt Betriebssystem-Funktionalitäten bereit. Wird hier genutzt, um durch Verzeichnisse zu navigieren
- Filetype: Wird zur Dateityp-Erkennung genutzt [10]
- RE: Das re-Modul stellt Funktionen bereit, mit denen mit regulären Ausdrücken gearbeitet werden kann. In der Beispiel-Implementierung werden mit diesem Modul die Tokens ermittelt. [python_re]
- Platform: Prüfung, auf welchem Betriebssystem das Programm läuft, da Windows- und Linux-Systeme verschiedene Dateisysteme nutzen
- Operator: Dieses Modul exportiert effiziente Funktionen, die den eigentlichen Operatoren von Python entsprechen. Es findet nur Anwendung in der Sortierung der zurückgebenden Dokumente in der retrieve-Methode.
- NLTK: Stellt Funktionen zur Verfügung, die es ermöglichen mit menschlichen Sprachdaten zu arbeiten.

Die Module Apache Tika, Filetype sowie NLTK müssen per pip installiert werden und sind damit nicht standardmäßig in Python integriert.

5.2.1 Apache Tika

Bei Apache Tika handelt es sich um eine Bibliothek, um Inhalte aus Dateien zu erkennen und zu analysieren. Es ist in der Lage Text und Metadaten aus verschiedenen Arten von Dateien zu extrahieren. Apache Tika liefert einen Parser, mit dessen Hilfe der Text aus - unter anderem - PDF-Dateien extrahiert werden kann. Mit dem Aufruf *parser.from_file(file)* kann ein PDF-Dokument in Text umgewandelt werden.

```
from tika import parser
import filetype
import math
import os
import string
import platform
import operator
import re
from nltk.tokenize import RegexpTokenizer
```

Abbildung 5.1: Imports

5.2.2 filetype

Mittels filetype ist es möglich, unabhängig von der Dateieindung, den Typ einer Datei zu ermitteln. Dies hat den Vorteil, dass die Suchmaschine sowohl unter Windows, als auch unter Unix-Systemen, alle PDF-Dateien finden kann, da unter Unix die Dateieindung keine garantierten Rückschlüsse auf den Typ der Datei zulässt.

5.2.3 NLTK

NLTK (natural language toolkit) ist eine Bibliothek für Python, die für die Verarbeitung natürlicher Sprachen eingesetzt wird. Dabei bietet die Bibliothek Funktionen für unter anderem Textklassifikation, Tokenisierung und Stemming. In diesem Beispiel wird NLTK verwendet, um die Eingabetexte der Dokumente und die Eingaben des Nutzers zu Tokens zu verarbeiten. Dazu wird eine Klasse verwendet, die auf Basis von Regular Expressions arbeitet. Mehr dazu wird anhand der Implementierung gezeigt.

5.2.4 Platform

Das Platform-Modul wird in dieser Implementierung dazu verwendet, festzustellen, auf welchem Betriebssystem die lokale Suchmaschine ausgeführt wird. Dies muss ermittelt werden, da auf Windows und Linux unterschiedliche Dateisysteme arbeiten. In Linux ist das Startverzeichnis immer das Root-Verzeichnis („/“). Unter Windows gibt meist mehrere Partitionen, die alle durchsucht werden müssen.

```
class Document:

    def __init__(self, url, length, docId, termList):
        self.url = url
        self.length = length
        self.docId = docId
        self.score = 0.
        self.termList = termList
```

Abbildung 5.2: Dokumentenklasse

5.3 Die Document-Klasse

Das Speichern der für das Retrieval wichtigen Informationen, geschieht mittels einer Document-Klasse (siehe Abbildung 5.2). Diese Klasse hält unter anderem alle Member-Variablen, die wichtig sind, um das TF-IDF-Maß berechnen zu können. Die Member-Variablen der Klasse sind:

- *url*: String mit dem Pfad zum Dokument, welches von dieser Instanz repräsentiert wird
- *length*: Integer, welcher die Anzahl der Wörter in diesem Dokument darstellt
- *docId*: Integer, welcher dem Dokument einen eindeutigen Identifikator gibt
- *score*: Float, welcher die Gewichtung der Dokument-Instanz zu einer Anfrage angibt
- *termList*: Eine Liste von Strings, die die Terme des Dokumentes entsprechen. Diese wird durch die Methode *_preprocess* erzeugt (siehe Abschnitt 5.4.2)

Die Member-Variable *url* soll am Ende des Retrievalprozesses zurückgegeben werden, da der Nutzer durch den Pfad direkten Zugriff auf das Dokument bekommt. Die Member-Variablen *length*, *termList* und *score* werden für die Berechnung des TF-IDF-Maßes benötigt und so auch für den Retrievalprozess. Die Member-Variable *docId* wird in der Indexklasse einerseits genutzt, um die Identifier in dem *invIndex*-Dictionary den jeweiligen Termen zuzuordnen. Des Weiteren erfolgt die Zuordnung der jeweiligen ID zu ihrer Document-Instanz über die *docHashmap* in der Indexklasse (siehe *invIndex* und *docHashmap* in Abschnitt 5.4).

```
def tf_idf(self, queryTerms, df, fileCount):  
  
    tfDict = {}  
    for term in queryTerms:  
        tfDict[term] = 0  
  
    for term in self.termList:  
        if term in queryTerms:  
            tfDict[term] += 1  
  
    for key, value in df.items():  
        idf = math.log(fileCount/value+1)  
        tfDict[key] *= idf  
  
    self.score = sum(tfDict.values())  
  
Document.tf_idf = tf_idf
```

Abbildung 5.3: TF-IDF-Methode

5.3.1 TF-IDF

Die Document-Klasse hält neben den benötigten Attributen auch die Scoringmethode *tf_idf* (siehe Abbildung 5.3). Dies ist die Implementierung des TF-IDF-Maßes und berechnet für jedes Dokument die Gewichtung für eine gegebene Suchquery. Hierbei benötigt die Funktion die Terme der Suchquery, welche über das Attribut *queryTerms* übergeben werden. Das Attribut *df* steht für die Document Frequency, welches für jeden Term die Anzahl der gefunden Dokumente in Form eines Dictionaries beinhaltet. Über das Attribut *fileCount* wird die Anzahl der Dokumente in der Kollektion, hier die Dateien die im invertierten Index aufgenommen wurden, übergeben.

Als erstes wird das Dictionary *tfDict* für die Term Frequency erstellt und mit den Termen der Query als Schlüssel und den Werten 0 initialisiert. In der nächsten For-Schleife wird über *self.termList* iteriert, welches die Liste der Terme des Dokumentes ist. Wenn ein Term auch in der Suchquery enthalten ist, also *term in queryTerms* zu true evaluiert wird, wird im *tfDict* der Wert des gefunden Terms um eins aufaddiert. Am Ende enthält das *tfDict* die Terme der Suchquery als Schlüssel und die Anzahl ihrer Vorkommnisse im Dokument als Wert, also die Term Frequency. Als letztes muss jede Term Frequency mit der Inverse Document Frequency multipliziert werden. Dafür wird über das Dictionary *df*

```
class Index:

    def __init__(self):
        self.invIndex = {}
        self.fileCount = 0
        self.docHashmap = {}
```

Abbildung 5.4: Indexklasse

iteriert und für jedes Key-Value-Paar die Inverse Document Frequency ausgerechnet und auf die jeweilige Term Frequency multipliziert. Am Ende wird die Summe aller TF-IDF-Maße der *score*-Variable zugeordnet und bilden so die finale Gewichtung für die gegebene Suche.

5.4 Der Index

Die Index-Klasse beinhaltet alle Methoden um den invertierten Index aufzubauen und die *retrieve*-Methode, welche die Dokumente, sortiert nach deren Gewicht, zu einer gegebenen Query zurückgibt (siehe Abbildung 5.4). Folgenden Member-Variablen werden für den invertierten Indexaufbau und die *retrieve*-Methode benötigt:

- *invIndex*: Ist ein Dictionary mit einem Term als Schlüssel und einer Menge von Dokumenten-IDs als Wert
- *fileCount*: Zählt beim Aufbau des invertierten Indexes die Dokumente, die in diesem aufgenommen wurden
- *docHashmap*: Ist ein Dictionary welches eine Dokumenten-ID ihrer zugehörige Document-Klasseninstanz zuordnet

Die Member-Variable *invIndex* ist der invertierte Index welcher bei der *retrieve*-Methode die Dokumenten-IDs zu einen gegebenen Term aus der Suchquery zurückgibt. Des Weiteren wird die Member-Variable *docHashmap* dafür genutzt, um über die Dokumenten-IDs auf die jeweiligen Dokumentinstanzen zuzugreifen und so auf ihren Inhalt und die Scoringfunktion *tf_idf*. Die Member-Variable *fileCount* gibt die Größe der Kollektion wieder und wird bei der Berechnung des TF-IDF-Maßes benötigt.

5.4.1 buildIndex

Die Methode *buildIndex* baut den invertierten Index auf und nutzt dafür die *__getStartDirectories*- und die *__addToIndex*-Methode. Dabei kann der *buildIndex*-Methode optional Startverzeichnisse in Form einer Liste mit raw Strings mitgegeben werden. In der ersten If-Abfrage wird geprüft, ob die Liste *startDirectories* leer ist. Dies ist der Fall, wenn der Nutzer keine Startverzeichnisse mitgegeben hat. Dann werden die Startverzeichnisse mithilfe der *__getStartDirectories*-Methode ermittelt.

Der erste Schritt stellt das Iterieren über alle Startverzeichnisse der Liste *startDirectories* dar. Anschließend werden für das Startverzeichnis, und alle darunterliegenden Verzeichnisse, bis zur untersten Ebene, die Dateien mithilfe der *os.walk*-Funktion geholt. Für jede Datei wird dann mit den Funktionen *os.path.abspath* und *os.path.join* der absolute Pfad gebildet. Danach wird durch das Modul *filetype* ermittelt, ob es sich um ein PDF-Dokument handelt. Ist ein Dokument vom Typ PDF, wird mithilfe des *Tika*-Moduls, genauer gesagt mit dem *parser*-Objekt, der Text aus dem PDF-Dokument extrahiert. Dies geschieht in dem die Methode *from_file* des Parsers, mit dem Pfad der Datei als Argument, aufgerufen wird. Bei der Rückgabe kann mithilfe des Schlüssels *content* auf den Text zugegriffen werden, welcher dann in der Variable *rawText* gespeichert wird.

Für jede entdeckte PDF-Datei wird die Member-Variable *fileCount* um eins erhöht. Da die Zahl sich bei jedem gefunden Dokument ändert, wird diese gleich als Dokumenten-ID genutzt, da sie für jedes Dokument eindeutig ist. Im Folgenden Schritt wird der *rawText* mithilfe der *__preprocessText*-Methode normalisiert, in eine Liste von Tokens geteilt und in der Liste *processedText* gespeichert. Nun sind alle Daten vorhanden um eine Document-Instanz zu erstellen. Diese enthält, wie in Abschnitt 5.3 schon beschrieben, den Pfad zum Dokument, die Anzahl der Wörter, die Dokumenten-ID und den Text als Liste von Tokens.

Im vorletzten Schritt wird die erstellte Document-Instanz seiner Dokumenten-ID in der Member-Variable *docHashmap* zugeordnet, um später auf diese Instanz zurückgreifen zu können. Als letztes wird der invertierte Index, mithilfe der *__addToIndex*-Methode, der Tokenliste *processedText* und der Dokumenten-ID, aktualisiert.

```
def buildIndex(self, startDirectories=[]):

    if not startDirectories:
        startDirectories = self._getStartDirectories()

    for directory in startDirectories:
        for root, _, files in os.walk(directory):
            for file in files:

                path = os.path.abspath(os.path.join(root, file))

                try:
                    if filetype.guess(path).mime == "application/pdf":

                        fileData = parser.from_file(path)
                        rawText = fileData['content']
                        self.fileCount += 1

                        processedText = self._preprocessText(rawText)
                        document = Document(path, len(processedText),
                                           self.fileCount, processedText)
                        self.docHashmap[self.fileCount] = document
                        self._addToIndex(self.fileCount, processedText)
                except:
                    continue

Index.buildIndex = buildIndex
```

Abbildung 5.5: BuildIndex-Methode

```
def _getStartDirectories(self):

    if platform.system() == "Linux":
        directories = ["/"]

    elif platform.system() == "Darwin":
        directories = ["/"]

    elif platform.system() == "Windows":
        directories = ['%s:\\' % d for d in string.ascii_uppercase
                       if os.path.exists('%s:' % d)]

    else:
        raise EnvironmentError

    return directories

Index._getStartDirectories = _getStartDirectories
```

Abbildung 5.6: GetStartDirectories-Methode

5.4.2 Hilfsmethoden

In diesem Abschnitt werden die drei Hilfsmethoden `_getStartDirectories`, `_preprocessText` und `_addToIndex` vorgestellt, welche in der `buildIndex`-Methode genutzt werden.

`_getStartDirectories`

Die Methode `_getStartDirectories` liefert eine Liste der Start-Verzeichnisse, abhängig vom Betriebssystem auf dem die Suchmaschine läuft. In diesen Verzeichnissen wird dann in der `buildIndex`-Methode rekursiv nach PDF-Dateien gesucht, welche in den Index mit einfließen. Falls das zugrunde liegende Betriebssystem ein Linux-basiertes oder Mac-basiertes System ist, wird die Liste `[, "/"]` zurückgegeben, da das Verzeichnis `/` immer das Root-Verzeichnis ist. Bei einem auf Windows basierenden Systemen gibt es wiederum mehrere Partitionen, welche immer mit einem Großbuchstaben abgekürzt werden. Dementsprechend gibt es unter Windows auch mehrere Start-Verzeichnisse.

Zuerst wird mithilfe der `system`-Funktion des `platform`-Moduls geprüft welches Betriebssystem vorliegt. Bei einem auf Linux- (Rückgabe „Linux“) oder Mac-basierenden System (Rückgabe „Darwin“) wird einfach eine List mit dem Element `/` erstellt.

```
def _addToIndex(self, documentID, terms):  
  
    for term in terms:  
  
        try:  
            docSet = self.invIndex[term]  
            docSet.add(documentID)  
            self.invIndex.update({term : docSet})  
  
        except KeyError:  
            self.invIndex[term] = {documentID}  
  
Index._addToIndex = _addToIndex
```

Abbildung 5.7: AddToIndex-Methode

Bei einem auf Windows basierenden Betriebssystem (Rückgabe „Windows“) ist das Erstellen der Startverzeichnisse aufwändiger. Hierbei werden alle Großbuchstaben, im Code durch *string.ascii_uppercase* aufrufbar, darauf geprüft eine Partition zu sein. Dies wird mithilfe der *os.path.exists*-Methode realisiert. Ist ein Großbuchstabe tatsächlich eine Partition auf dem Computer, so wird er in der Liste gespeichert. Jedoch wird an den Großbuchstaben noch der String „:\“ angehängen, damit die *buildIndex*-Methode mit den Elementen der Liste als Start-Verzeichnisse arbeiten kann.

`_addToIndex`

Die Methode *_addToIndex* fügt dem Index für jeden Term *t* die Dokumenten-ID an, in welchem der Term *t* gefunden wurde. Hierfür bekommt die Methode eine Liste von Termen die in einem PDF-Dokument vorkommen über den Parameter *terms* und die zum Dokumente gehörige Dokumenten-ID als Parameter *documentID* übergeben.

Für jeden Term in der Liste *terms* wird dazu der zum Term gehörige Eintrag im invertierten Index nachgeschlagen. Schlägt der Versuch fehl, da noch kein Eintrag des Terms in dem invertierten Index vorhanden ist, wird ein neuer Eintrag für diesen Term und der übergebenen Dokumenten-ID erstellt. Wird ein Eintrag gefunden, wird die übergebende Dokumenten-ID der schon vorhandene Menge hinzugefügt und der invertierte Index wird geupdatet.

`_preprocessText`

Diese Methode dient der Vorverarbeitung der Texte, die in den PDF-Dokumenten stehen, wird aber auch für die Suchquery genutzt. Hierfür wird der Text eines PDF-Dokumentes an den Parameter *text* übergeben. Als erster Schritt wird der gesamte Text in Lower-Case (Kleinschreibung) gesetzt, damit später bei der Suche die Groß- bzw. Kleinschreibung irrelevant ist. Das Ergebnis wird in der Variable *lowerText* gespeichert. Im nächsten Schritt werden alle Zahlen aus *lowerText* entfernt und in *prepText* gespeichert, da Zahlen für die Textsuche nicht von Bedeutung sind.

Als nächstes wird mithilfe der Klasse *RegexTokenizer*, die durch die NLTK-Bibliothek zur Verfügung gestellt wird, der String *prepText* in eine Liste von Tokens aufgespalten. Was als Token gewertet wird, wird mithilfe einer *regular expression* definiert, im Deutschen regulärer Ausdruck genannt. Ein regulärer Ausdruck ist eine Zeichenkette, welche eine Menge von bestimmten Zeichenketten beschreibt. Der gewünschte reguläre Ausdruck wird dem Konstruktor der *RegexTokenizer*-Klasse in Form eines raw Strings übergeben. Ein raw String ist ein String, welcher mit einem *r* am Anfang gekennzeichnet ist und ein Backslash als ein Literal und kein Escape-Zeichen behandelt. Dies ist bei regulären Ausdrücken nützlich, da in diesen verstärkt mit Backslashes gearbeitet wird.

Im Folgenden wird der raw String bzw. reguläre Ausdruck näher betrachtet, um zu verstehen, was als ein Token gewertet wird. Der erste Teil des regulären Ausdrucks `[a-zA-Z]+` definiert alle Buchstabenketten mit einen oder mehreren Elementen, die mit einem Bindestrich bei einem Zeilenumbruch enden, als Token. Die eckigen Klammern werden bei regulären Ausdrücken genutzt, um eine Zeichenauswahl zu definieren. Das bedeutet, dass ein Zeichen aus dieser Auswahl dann an dieser Stelle steht. Mithilfe von Quantoren kann definiert werden, wie viele Zeichen einer Auswahl hintereinander stehen dürfen. Das Pluszeichen ist genau so ein Quantor, welcher aussagt, dass mindestens ein oder mehrere Zeichen der Zeichenauswahl hintereinander vorkommen muss. Die Funktion des Dollarzeichens hängt von einem Flag, dem Multiline-Flag ab. Dieses ist in NLTK standardmäßig gesetzt und bewirkt, dass das Dollarzeichen für das Ende einer Zeile steht.^[9] Sonst würde das Dollarzeichen das Ende eines Wortes markieren. Dadurch das ein Bindestrich vor das Dollarzeichen des regulären Ausdrucks gesetzt haben, bedeutet der reguläre Teilausdruck

-\$, dass die letzte Zeichenkette einer Zeile auf einem Bindestrich endet.

Bei dem `|`-Zeichen handelt es sich um eine logische Oder-Verknüpfung, die es ermöglicht, mehrere reguläre Ausdrücke zu verknüpfen. In diesem Fall `\w+`. Der zweite reguläre Ausdruck `\w+` definiert alle alphanumerischen Zeichenketten mit einem oder mehreren Elementen als Token. Hierbei ist `\w` eine vordefinierte Zeichenklasse für den regulären Ausdruck `[a-zA-Z_0-9]` und beinhaltet außer alphanumerische Werte auch noch den Unterstrich. Das Pluszeichen ist hier wieder der Quantor, welcher aussagt, dass aus dieser Zeichenklasse ein oder mehrere Zeichen hintereinander vorkommen muss. Mithilfe der *tokenize*-Methode wird der reguläre Ausdruck auf den String *prepText* angewendet. Jeder Substring des mitgegebenen Strings, der den regulären Ausdruck erfüllt, wird an die Liste *tokenList* angefügt.

Der Grund warum die Wörter, die auf einem Bindestrich enden, bei der Tokenerzeugung extra beachtet werden, ist der, dass die Wörter, welche bei Zeilenumbrüchen getrennt werden, wieder zusammengefügt werden sollen. In der for-Schleife werden diese Tokens auf die Eigenschaft hin, auf einem Bindestrich zu enden, geprüft und gegebenenfalls zusammengesetzt. Dazu wird der Bindestrich aus dem Token entfernt und mit dem nächsten Token in der Tokenliste verknüpft (`token[:-1]+tokenList[index+1]`). Die Tokenliste wird darauf hin aktualisiert. Der neu zusammengesetzte Token ersetzt den Token, welcher den Bindestrich enthielt, und der darauffolgende Token der Liste wird gelöscht, da er schon mit dem Token davor zusammengesetzt wurde.

Diese Methode, die Wörter mit einem Bindestrich am Zeilenende mit dem nächsten Token zusammenzufügen, ist jedoch nicht immer korrekt. Es kann auch folgender Fall eintreten: Eine Zeile endet zum Beispiel mit *Damen-* und die nächste Zeile geht mit *und Herrenschuhe* weiter. In diesem Fall ist der Bindestrich gewollt, der Algorithmus fügt jedoch die Wörter *Damen* und *und* zu einem Wort zusammen. Da davon auszugehen ist, dass dieser Fall selten eintritt, wurde er vernachlässigt.

```
def _preprocessText(self, text):

    lowerText = text.lower()

    prepText = re.sub(r'\d+', '', lowerText)

    tokenizer = RegexpTokenizer(r'[a-zA-Z]+-$|\w+')
    tokenList = tokenizer.tokenize(prepText)

    for token in tokenList:
        if token[-1] == '-':

            index = tokenList.index(token)
            compositeWord = token[:-1]+tokenList[index+1]
            tokenList[index] = compositeWord
            del tokenList[index+1]

    return tokenList

Index._preprocessText = _preprocessText
```

Abbildung 5.8: PreprocessText-Methode

5.4.3 retrieve

Die *retrieve*-Methode dient der Suche nach Dokumenten anhand einer eingegebenen Suchanfrage und stellt die Schnittstelle zum Nutzer dar. Die Idee dabei ist, dass der Nutzer ein oder mehrere Schlagworte, in Form eines Strings, der Methode übergeben kann. Auf deren Basis werden die gefundenen Dokumente, nach Relevanz sortiert, zurückgeliefert.

Zunächst wird der Suchstring des Nutzers mittels der bereits bekannten Methode *__preprocessText* normalisiert und in eine Liste von Termen zerteilt. Das Ergebnis wird in der Variable *processedStrings* abgespeichert. Im zweiten Schritt werden alle genutzten Variablen deklariert. Mithilfe der For-Schleife wird über die Liste von Termen *processedStrings* iteriert, um über den invertierten Index für jeden Term die Menge der Dokumenten-IDs zu beschaffen. Des Weiteren wird für jeden Term die Länge seiner Menge von Dokumenten-IDs im Dictionary *df* abgespeichert und repräsentiert so die Document Frequency. Die Variable *result* stellt die Menge da, die alle zur Suchquery gefundenen Dokumenten-IDs enthält. So wird diese Variable mit jeder Menge von Dokumenten-IDs vereinigt. Da es sich hier um eine Menge handelt kann ausgeschlossen werden, dass Dokumenten-IDs doppelt vorkommen. Existiert der Term *word* nicht als Key im Index, da der Nutzer ein Wort

eingetragen hat, welches kein Dokument beinhaltet, wird beim nächsten Term der Liste *processedStrings* fortgefahren.

In der nächsten For-Schleife wird für jede Dokumenten-ID in der Liste *result*, mithilfe der *docHashmap*, die passende Documentinstanz geholt, um für jedes Dokument das TF-IDF-Maß auszurechnen. Wurde dies mithilfe der *tf_idf*-Methode getan, wird das Gewicht zu der jeweiligen Dokumenten-ID in dem Dictionary *weightedDocs* gespeichert. Dieses wird mit der *sorted*-Funktion von Python nach den Gewichten sortiert. Die Sortierung ist jedoch aufsteigend, obwohl es gewünscht ist, die Dokumente mit einem höheren Gewicht oben zu haben. Dafür wird die Variable *reverse* auf True gesetzt, damit die Einträge absteigend sortiert werden.

In der letzten For-Schleife werden die absoluten Pfade (*ind.docHashmap[Key].url*), in der sortierten Reihenfolge, und die Gewichte der Dokumente für die Query in die Liste *resultList* gespeichert. Diese wird am Ende zurückgegeben.

5.5 Ausführung der Suchmaschine

In diesem Abschnitt soll kurz erläutert werden, wie die oben erstellten Klassen genutzt werden. Als erstes muss eine neue Instanz der Indexklasse erstellt werden. Über diese Instanz wird dann die *buildIndex*-Methode aufgerufen. Optional kann hier eine Liste der gewünschten Startverzeichnisse als raw Strings übergeben werden. Wenn der invertierte Index aufgebaut wurde, kann über die Indexinstanz die *retrieve*-Methode mit der gewünschten Suchquery aufgerufen werden. Diese liefert eine nach Gewichten sortierte Liste, mit den Pfaden zu den gefundenen Dokumenten. Wenn diese nicht leer ist, können mithilfe einer For-Schleife die Elemente ausgegeben werden.

```
def retrieve(self, searchString):

    processedStrings = self._preprocessText(searchString)
    result = set()
    df, weightedDocs = {}, {}
    resultList = []

    for word in processedStrings:
        try:
            documents = self.invIndex[word]
            df[word] = len(documents)
            result = result.union(documents)
        except KeyError:
            continue

    for document in result:
        doc = ind.docHashmap[document]
        doc.tf_idf(processedStrings, df, self.fileCount)
        weightedDocs[doc.docId] = doc.score

    sortedDocs = sorted(weightedDocs.items(), key=operator.itemgetter(1),
                        reverse=True)

    for key, value in sortedDocs:
        resultList.append([ind.docHashmap[key].url, value])

    return resultList

Index.retrieve = retrieve
```

Abbildung 5.9: Retrieve-Methode

```
ind = Index()
ind.buildIndex([r"Opt\Start\1", r"Opt\Start\2", ...])

resultSet = ind.retrieve("This is a test")
if resultSet:
    print('Pfad, \t Gewicht')
    for url, score in resultSet:
        print(url, '\t', score)
```

Abbildung 5.10: Execute

6 Verbesserungen

In der Vorgestellten Implementierung gibt es einige Probleme und daraus folgende Verbesserungen, die in diesem Abschnitt besprochen werden.

6.1 Probleme der aktuellen Beispiel-Implementierung

Dieser Abschnitt wird eine kurze Problemanalyse der aktuellen Implementierung der lokalen Suchmaschine liefern. Zu diesem Zweck soll die NPL-Test-Collection herangezogen werden. Mit deren Hilfe kann ermittelt werden, wie „gut“ die Suchmaschine ist. Dazu werden zwei Messungen vorgenommen:

- Die *precision* gibt den Anteil an, wie viele der gefundenen Dokumente relevant sind.
- Der *recall* gibt an, wie viele der relevanten Dokumente gefunden werden.

6.1.1 Analyse der Precision

Die Precision ist wie folgt definiert:

$$precision = \frac{r_1}{r_2},$$

wobei r_1 die Anzahl der ermittelten, relevanten Dokumente ist und r_2 die Anzahl aller ermittelten Dokumente.

In der NPL-Collection sind 11.429 Dokumente vorhanden, zudem sind 93 Queries mitgeliefert. Den Queries liegen darüber hinaus alle Dokumenten-IDs bei, die für die Query relevant sind. Falls die Suchmaschine genau die einer Query zugeordneten Dokumenten

zurückliefert, wird die *precision* 1. Liefert der Bruch $\frac{r_1}{r_2} \geq 0,8$, ist das Ergebnis bereits ausreichend gut.

Test und Auswertung der Ergebnisse

Zunächst werden 10 Queries getestet und auf deren Basis wird ein Wert für die Precision gemessen. Der durchschnittliche Wert der Precision wird mit $\frac{1}{10} * \sum_1^{10} \frac{r_{1_i}}{r_{2_i}}$ berechnet.

Der Wert der Precision ist $< 0,1$. Dieser Wert ist sehr niedrig. Dies hat den Grund, dass die Suchmaschine für jedes Wort der Query die Menge der Dokumente zurückliefert, in der das Wort steht. Durch den Union-Operator werden die Ergebnis-Mengen der einzelnen Wörter lediglich vereinigt. Dadurch reicht es aus, dass ein Wort der Query in nur einem Dokument vorkommt, um in der Ergebnis-Menge zu stehen.

Lösungsvorschläge

Das Bilden von N-Grammen ist ein Weg dieses Problem anzugehen. Um das Konzept des N-Gramms zu verstehen, wird ein Beispiel gegeben:

6.1.2 Analyse des Recalls

Der Recall ist definiert als

$$recall = \frac{r_1}{r_2},$$

wobei r_1 die Anzahl der ermittelten, relevanten Dokumente ist und r_2 die Anzahl relevanter Dokumente ist. Der *recall* wird 1, wenn die Menge der zurückgegebenen Dokumente genau der Menge der relevanten Dokumenten in der Collection entspricht. Gilt $recall \geq 0,8$, ist der Wert ausreichend gut.

Test und Auswertung der Ergebnisse

Auch hier 10 Queries getestet und auf deren Basis wird ein Wert für den Recall gemessen. Der durchschnittliche Wert des Recalls wird mit $\frac{1}{10} * \sum_1^{10} \frac{r_{1_i}}{r_{2_i}}$ berechnet.

Im Gegensatz zur Precision liefert der Recall einen Wert $\gg 0,95$. Diese Tatsache ist wenig überraschend, da - wie im letzten Abschnitt gesehen - nahezu alle Dokumente zurückgeliefert werden. Daher sind in den meisten Fällen alle Dokumente, die für die Query relevant sind, in der Ergebnis-Menge.

Da der Recall für sich gesehen einen sehr guten Wert liefert, muss hier keine Verbesserung vorgenommen werden.

6.2 Aufbaudauer des Index

Eine Verbesserung betrifft die Geschwindigkeit, in der der Index aufgebaut wird. Für 189 PDF-Dokumente benötigt die Beispiel-Implementierung ca 1 Minute und 17 Sekunden. Wird ein System mit ca. 105GB nach PDF-Dateien durchsucht und in den Index aufgenommen, werden auf einem System mit einem Intel Core i5 der 8. Generation und 8 Gigabyte RAM ca. 40 Minuten benötigt. Für einen sinnvollen Einsatz die Suchmaschine somit nicht geeignet. Zwei Ideen, diese Zeit zu senken, sind Parallelisierung und Persistierung.

6.2.1 Persistierung

Die Idee, die wahrscheinlich am wirksamsten ist, besteht darin, den Index nach oder bereits während des Aufbaus, zu persistieren. Durch die Persistierung kann bei einem Neustart der Suchmaschine eine Datei eingelesen werden, welche die Index-Daten beinhaltet. Dadurch ist es nicht nötig, beim Start durch alle Pfade zu gehen und jede Datei zu öffnen, um festzustellen, ob es sich um eine PDF-Datei handelt. Pro Öffnen einer Datei werden mindestens zwei Systemaufrufe benötigt (öffnen und schließen der Datei). Zudem werden weitere Systemaufrufe benötigt, wenn Verzeichnisse geöffnet werden, da diese in der Regel ebenfalls als spezielle Dateien behandelt werden. Das hat zur Folge, dass viele Kontextwechsel stattfinden und das Programm bei jedem Systemaufruf blockiert wird.

Dadurch, dass bei einem persistenten Index initial nur eine Datei geöffnet werden muss, kann ein Großteil der Systemaufrufe und Kontextwechsel gespart werden. Damit kann die

benötigte Zeit zum Aufbau des Index gesenkt werden.

Ein Problem, was hier gelöst werden muss ist, alle Dokumente ausfindig zu machen, die noch nicht im Index aufgenommen wurden. Für dieses Problem kann Parallelisierung genutzt werden, die im Folgenden Abschnitt behandelt wird.

6.2.2 Parallelisierung

Durch Nutzung von Parallelisierung kann der Aufbau des Index ebenfalls beschleunigt werden. Da die Beispiel-Implementierung auf Python basiert, ist jedoch zu beachten, dass Python-Threads keine Vorteile aus Multi-Prozessorsystemen ziehen kann. Das bedeutet, dass mehrere Python-Threads nicht parallel auf verschiedenen Kernen laufen können. Dennoch bringt Parallelisierung einen Geschwindigkeits-Vorteil.

Wird ein Systemaufruf ausgelöst, um beispielsweise eine Datei zu öffnen, wird, sofern die Suchmaschine Single-Threaded ist, das komplette Programm blockiert. Werden jedoch Threads genutzt, kann nur der Thread blockiert werden, der den Systemaufruf verursacht hat. An dessen Stelle kann dann ein weiterer Thread treten, der in der Zwischenzeit rechnen kann. Dadurch wird weniger Zeit mit Warten verschwendet und der Index-Aufbau wird beschleunigt.

Ist der Index persistent gespeichert, ist es möglich, in regelmäßigen einen oder mehrere Threads zu starten und somit „im Hintergrund“ nach neuen Dateien gesucht werden, die noch nicht im Index vorhanden sind.

Durch die Kombination beider Ideen, sollte der Aufbau des Index spürbar beschleunigt werden.

Literatur

- [1] Nagy Istvan. *Indexierung mittels invertierter Dateien*. 2001. URL: <http://www.cis.uni-muenchen.de/people/Schulz/SeminarSoSe2001IR/Nagy/node4.html>.
- [2] Ingo Frommholz. *Information Retrieval*. 2007. URL: http://www.is.informatik.uni-duisburg.de/courses/ie_ss07/folien/folien-ir.pdf.
- [3] Andreas Henrich. *Information Retrieval 1*. 2008.
- [4] Hinrich Schütze Christopher D. Manning Prabhakar Raghavan. *Introduction to Information Retrieval*. 2009.
- [5] Vaidehi Joshi. *Trying to Understand Tries*. 2017. URL: <https://medium.com/basecs/trying-to-understand-tries-3ec6bede0014>.
- [6] *Vektorraum-Retrieval*. 2017. URL: <https://de.wikipedia.org/wiki/Vektorraum-Retrieval>.
- [7] o.A. *Trie*. 2018. URL: <https://de.wikipedia.org/wiki/Trie>.
- [8] o.V. *Apache Tika - a content analysis toolkit*. 22. Dez. 2018. URL: <https://tika.apache.org/>.
- [9] o.V. *source code for nltk.tokenize.regexp*. 17. Nov. 2018. URL: https://www.nltk.org/_modules/nltk/tokenize/regexp.html.
- [10] o.V. *fieltype 1.0.5*. 1. März 2019. URL: <https://pypi.org/project/filetype/>.
- [11] freeCodeCamp. *Introduction to Trie*. URL: <https://guide.freecodecamp.org/miscellaneous/data-structure-trie/>.
- [12] Karin Haendelt. *Klassische Information Retrieval Modelle Einführung*.
- [13] o.V. *Information Retrieval*. URL: https://de.wikipedia.org/wiki/Information_Retrieval.