

Konzeption und Realisierung eines Systems zur Informationssuche in einem Dokumentenarchiv basierend auf Textinhalt und Metadaten.

Conception and Realization of an Information Retrieval System for a Document Archive based on Text Content and Metadata

Annika Kremer

Bachelor-Abschlussarbeit

Betreuer: Prof. Dr. Karl Hans Bläsius

Trier, Abgabedatum

Vorwort

Ein Vorwort ist nicht unbedingt nötig. Falls Sie ein Vorwort schreiben, so ist dies der Platz, um z.B. die Firma vorzustellen, in der diese Arbeit entstanden ist, oder einigen Leuten zu danken, die in irgendeiner Form positiv zur Entstehung dieser Arbeit beigetragen haben. Auf keinen Fall sollten Sie im Vorwort die Aufgabenstellung näher erläutern oder vertieft auf technische Sachverhalte eingehen.

Kurzfassung

In der Kurzfassung soll in kurzer und prägnanter Weise der wesentliche Inhalt der Arbeit beschrieben werden. Dazu zählen vor allem eine kurze Aufgabenbeschreibung, der Lösungsansatz sowie die wesentlichen Ergebnisse der Arbeit. Ein häufiger Fehler für die Kurzfassung ist, dass lediglich die Aufgabenbeschreibung (d.h. das Problem) in Kurzform vorgelegt wird. Die Kurzfassung soll aber die gesamte Arbeit widerspiegeln. Deshalb sind vor allem die erzielten Ergebnisse darzustellen. Die Kurzfassung soll etwa eine halbe bis ganze DIN-A4-Seite umfassen.

Hinweis: Schreiben Sie die Kurzfassung am Ende der Arbeit, denn eventuell ist Ihnen beim Schreiben erst vollends klar geworden, was das Wesentliche der Arbeit ist bzw. welche Schwerpunkte Sie bei der Arbeit gesetzt haben. Andernfalls laufen Sie Gefahr, dass die Kurzfassung nicht zum Rest der Arbeit passt.

The same in english.

Inhaltsverzeichnis

1	Einleitung und Problemstellung	1
1.1	Einleitung	1
1.2	Problemstellung	2
1.3	Teilprobleme	2
1.3.1	Unterscheidung zwischen Metadaten und Freitext	2
1.3.2	Metadatensuche	2
1.3.3	Freitextsuche	2
1.3.4	Verknüpfung mit UND/ODER	2
1.3.5	Benutzeroberfläche	2
2	Information Retrieval	3
3	Boolean Retrieval	4
3.0.1	Matrix	4
3.0.2	Posting List	4
3.0.3	Anfragenverknüpfung	4
4	Das Vektorraum-Modell	5
4.1	Funktionsprinzip	5
4.2	Term Frequency	5
4.3	Document Frequency	5
4.4	Inverted Document Frequency	5
4.5	$Tf \times idf$ Weighting	5
4.5.1	Formeln	5
4.6	Ähnlichkeitsfunktion	5
4.6.1	Euklidische Distanz	5
4.6.2	Cosine Similarity	5
4.6.3	Alternativen	5
5	Bewertung eines Information Retrieval Systems	6
5.1	Precision	6
5.2	Recall	6
6	Implementierung	7

7 Die Benutzeroberfläche	8
8 Zusammenfassung und Ausblick	9
Literaturverzeichnis	10
Glossar	11
Erklärung der Kandidatin / des Kandidaten	12

Abbildungsverzeichnis

Tabellenverzeichnis

Einleitung und Problemstellung

1.1 Einleitung

Tag für Tag treffen neue E-Mail Nachrichten mit den unterschiedlichsten Inhalten im Posteingang ein, weshalb die Anzahl gesendeter und empfangener Mails schnell unübersichtlich groß wird.

Will der Nutzer eine Weile später auf eine bestimmte Nachricht erneut zugreifen, steht er oftmals vor der Problematik, mit der manuellen Suche aufgrund der schier unübersichtlichen Datenmenge überfordert zu sein.

Besonders schwierig wird es für ihn, wenn er nach bestimmten Textinhalten sucht, sich jedoch weder an das genaue Datum noch an den Absender erinnert. An dieser Stelle hilft dem Nutzer ein Suchprogramm, mit dem er Suchkriterien beliebig logisch verknüpfen kann.

Weiß er beispielsweise noch, dass die Nachricht im Mai ankam und eine wichtige Adresse beinhaltet, könnte seine Suchanfrage „Datum = Mai und Freitext = Adresse“ lauten. Das Suchprogramm präsentiert dann eine Liste für ihn in Frage kommender Dokumente, was seine Suche auf einen kleinen Kreis relevanter Treffer einschränkt.

Hier rumlabern von wegen in E-Mails sucht man manchmal nach iwas, unübersichtlich, bla blubb

Begonnen werden soll mit einer Einleitung zum Thema, also Hintergrund und Ziel erläutert werden.

Weiterhin wird das vorliegende Problem diskutiert: Was ist zu lösen, warum ist es wichtig, dass man dieses Problem löst und welche Lösungsansätze gibt es bereits. Der Bezug auf vorhandene oder eben bisher fehlende Lösungen begründet auch die Intention und Bedeutung dieser Arbeit. Dies können allgemeine Gesichtspunkte sein: Man liefert einen Beitrag für ein generelles Problem oder man hat eine spezielle Systemumgebung oder ein spezielles Produkt (z.B. in einem Unternehmen), woraus sich dieses noch zu lösende Problem ergibt.

Im weiteren Verlauf wird die Problemstellung konkret dargestellt: Was ist spezifisch zu lösen? Welche Randbedingungen sind gegeben und was ist die Zielsetzung? Letztere soll das beschreiben, was man mit dieser Arbeit (mindestens) erreichen möchte.

1.2 Problemstellung

Teilweise strukturiert, teilweise Freitext, man will zeugs finden -; relativ allgemein halten

1.3 Teilprobleme

1.3.1 Unterscheidung zwischen Metadaten und Freitext

1.3.2 Metadatensuche

1.3.3 Freitextsuche

1.3.4 Verknüpfung mit UND/ODER

1.3.5 Benutzeroberfläche

Boolean Retrieval

3.0.1 Matrix

3.0.2 Posting List

3.0.3 Anfragenverknüpfung

Das Vektorraum-Modell

4.1 Funktionsprinzip

4.2 Term Frequency

4.3 Document Frequency

4.4 Inverted Document Frequency

4.5 $Tf \times idf$ Weighting

4.5.1 Formeln

4.6 Ähnlichkeitsfunktion

4.6.1 Euklidische Distanz

4.6.2 Cosine Similarity

4.6.3 Alternativen

Bewertung eines Information Retrieval Systems

5.1 Precisison

5.2 Recall

Implementierung

Die Benutzeroberfläche

Zusammenfassung und Ausblick

In diesem Kapitel soll die Arbeit noch einmal kurz zusammengefasst werden. Insbesondere sollen die wesentlichen Ergebnisse Ihrer Arbeit herausgehoben werden. Erfahrungen, die z.B. Benutzer mit der Mensch-Maschine-Schnittstelle gemacht haben oder Ergebnisse von Leistungsmessungen sollen an dieser Stelle präsentiert werden. Sie können in diesem Kapitel auch die Ergebnisse oder das Arbeitsumfeld Ihrer Arbeit kritisch bewerten. Wünschenswerte Erweiterungen sollen als Hinweise auf weiterführende Arbeiten erwähnt werden.

Literaturverzeichnis

- CDK02. COULOURIS, GEORGE, JEAN DOLLIMORE und TIM KINDBERG: *Verteilte Systeme: Konzepte und Design*. Addison-Wesley-Verlag, 2002.
- Che85. CHERITON, DAVID R.: *Preliminary Thoughts on Problem-oriented Shared Memory: A Decentralized Approach to Distributed Systems*, 1985.
- Mal97. MALTE, PETER: *Replikation in Mobil Computing*. Seminar No 31/1997, Institut für Telematik der Universität Karlsruhe, Karlsruhe, 1997.
<http://www.ubka.uni-karlsruhe.de/cgi-bin/psview?document=/ira/1997/31>.
- Mos93. MOSBERGER, DAVID: *Memory Consistency Models*. Technical Report 93/11, University of Arizona, November 1993.

A

Glossar

DisASter	DisASter (Distributed Algorithms Simulation Terrain), A platform for the Implementation of Distributed Algorithms
DSM	Distributed Shared Memory
AC	Linearisierbarkeit (atomic consistency)
SC	Sequentielle Konsistenz (sequential consistency)
WC	Schwache Konsistenz (weak consistency)
RC	Freigabekonsistenz (release consistency)

B

Erklärung der Kandidatin / des Kandidaten

☐ Die Arbeit habe ich selbstständig verfasst und keine anderen als die angegebenen Quellen- und Hilfsmittel verwendet.

☐ Die Arbeit wurde als Gruppenarbeit angefertigt. Meine eigene Leistung ist ...

Diesen Teil habe ich selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet.

Namen der Mitverfasser: ...

Datum

Unterschrift der Kandidatin / des Kandidaten