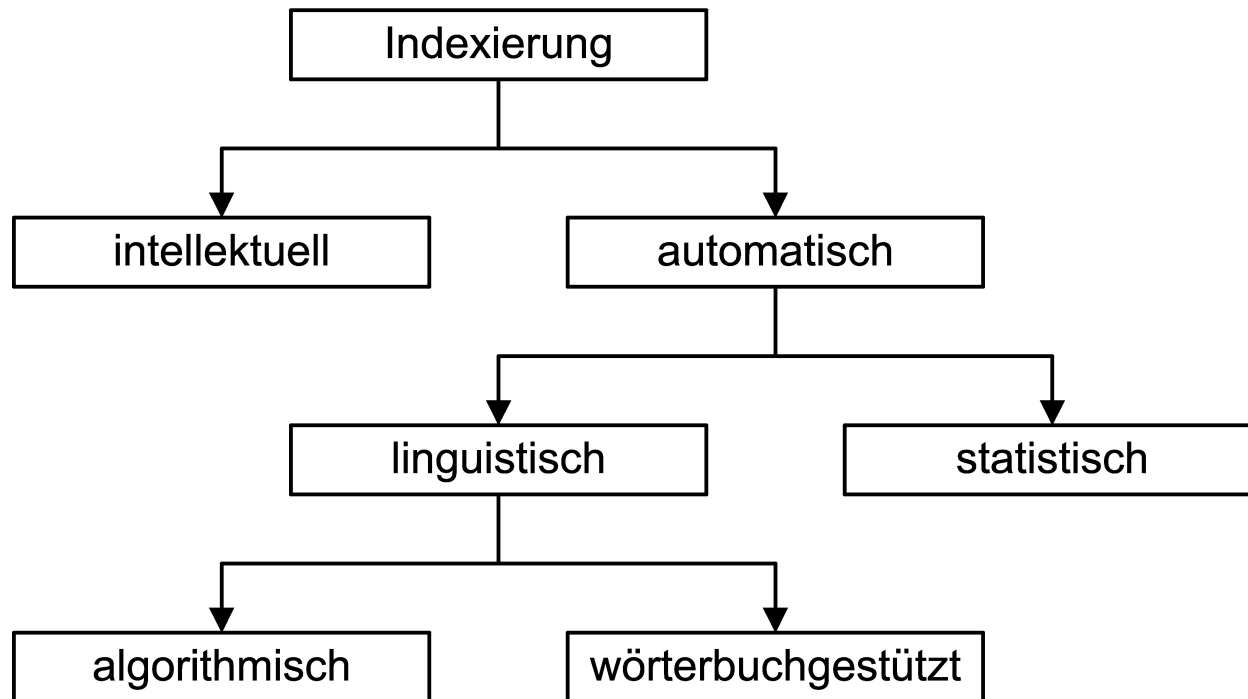


Automatisches Erschließen

Verfahren der Indexierung



Kapitel 5 - Abb. 5.1

Das Prinzip Index



Ein Satz, für uns bestehend aus Wörtern:

Für das IR-System besteht ein abgespeichertes Dokument nicht aus „Wörtern“ sondern lediglich aus einer Aneinanderreihung von Zahlencodes:



Speicherinhalt:

```
46C3BC72206461732049522D53797374656D20626573746568742065696E2061626765737065696368657274
657320446F6B756D656E74206E6174C3BC726C696368206E696368742061757320E2809E57C3B6727465726E
E2809C20736F6E6465726E206C656469676C696368206175732065696E657220416E65696E616E6465727265
6968756E6720766F6E205A61686C656E636F6465732E
```

Textmatching – Sequentielle Suche



Für das IR-System besteht ein abgespeichertes Dokument nicht aus „Wörtern“ sondern lediglich aus einer Aneinanderreihung von Zahlencodes:



→
46C3BC72206461732049522D53797374656D20626573746568742065696E2061626765737065696368657274
657320446F6B756D656E74206E6174C3BC726C696368206E696368742061757320E2809E57C3B6727465726E
E2809C20736F6E6465726E206C656469676C696368206175732065696E657220416E65696E616E6465727265
6968756E6720766F6E205A61686C656E636F6465732E

Textmatching – Indexsuche – invertierte Liste

Zeichenketten im Dokument

Für das IR-System besteht ein abgespeichertes Dokument nicht aus „Wörtern“ sondern lediglich aus einer Aneinanderreihung von Zahlencodes:



Index

TERM	DOK-NR
abgespeichertes	1,145,56398
Aneinanderreihung	1,3,189
aus	1,2,3,4,5,6...
besteht	1,8,15,18,89,...
das	1,2,3,4,5,6,7...
ein	1,2,3,4,5,6,7...
...	
Zahlencodes	1

allgemeine Form

Term	Dok-Nr	(Treffer)	(Position)
aus	1,2,3,...	54378	27,456,3...
...
term	1,...n	n	n

Automatische Indexerstellung

Zerlegung eines Textes in Zeichenketten

Vorgehensweise:

Definition von Zeichenkette, bspw. „Alles zwischen zwei Leerzeichen“

Extraktion aller Zeichenketten durch Algorithmisierung der Definition

Alphabetisches Sortieren des Ergebnisses

Tabelle 5.1 Manuelle erzeugte Indexterme

-	(open	Abstract:
algorithmische	an	Anwendungsgrenzen
auf.	Automatische	automatischen
automatischen	basierten	Bei
Beispielen	Beitrag	benannt.
beschrieben,	bzw.	das
den	der	der
Der	Der	des
deutschen	deutschsprachiger	die
Die	Dokumente	ein
ein	eine	einer
Einsatzmöglichkeiten	Einsatzszenarien	Entwicklung
erläutert:	Flexibilität	frei
Funktionalität	für	für
für	für	Grundformerkenung,
hohe	im	Indexierung
Indexierung	Indexierung	Information
ist	Kompositumerkennung	Kompositumzerlegung,
Konfigurierbarkeit	lexikalische	Lingo
Lingo	lingo	lingo
lingo	linguistisch	linguistische
Mehrwortgruppenerkennung,	mögliche	Nutzen
OCR-Fehlerkorrektur.	offene	open
Retrieval	Retrievalverbesserung	source
source)	Sprache.	standen
stehende	System	System
Systemaufbau	Systems	Titel:
und	und	und
und	unterschiedliche	verfügbares
Verfügung	von	von
von	Vordergrund.	vorgestellt
werden	wird	wird
Wortrelationierung,	zeigt	zur
zur		

Intellektuelle Registererstellung

Festlegung von Sucheinstiegen auf Textstellen zum Thema

Titel:

Lingo - ein **open source** System für die **Automatische Indexierung** deutschsprachiger Dokumente

Abstract:

Lingo ist ein frei verfügbares System (open source) zur automatischen Indexierung der deutschen Sprache. Bei der Entwicklung von lingo standen hohe Konfigurierbarkeit und Flexibilität des Systems für unterschiedliche Einsatzmöglichkeiten im Vordergrund. Der Beitrag zeigt den Nutzen einer **linguistisch basierten automatischen Indexierung** für das **Information Retrieval** auf. Die für eine Retrievalverbesserung zur Verfügung stehende linguistische Funktionalität von lingo wird vorgestellt und an Beispielen erläutert: **Grundformerkennung**, Kompositumerkennung bzw. **Kompositumzerlegung**, Wortrelationierung, lexikalische und **algorithmische Mehrwortgruppenerkennung**, OCR-Fehlerkorrektur. Der offene Systemaufbau von lingo wird beschrieben, mögliche Einsatzszenarien und Anwendungsgrenzen werden benannt.

Registereingang

Automatische Indexierung
- Grundformerkennung
- Kompositumerkennung
- Kompositumzerlegung
- Mehrwortgruppenerkennung
- Wortrelationierung
Information Retrieval
Lingo
Mehrwortgruppenerkennung
- algorithmisch
- lexikalisch
Open source

Abb. 5.3 Beispiele für ausgewählte Registerbegriffe

Registerausgang

nn
nn
nn
nn
nn
nn
nn
nn
nn
nn
nn

Vorgehensweise

Auswahl von wichtigen Wörtern im Text

Verwenden einer Funktion zur Markierung von Index- bzw. Registerbegriffen

Erzeugen des Index bzw Registers - alphabetisch sortiert, ggf. mit Hierarchie

Abb. 5.4 Beispiel für ein generiertes Register

ausgewähltes Vokabular

kontrolliertes Vokabular

Stichwortregister

Automatische Indexierung

- Grundformerkennung
- Kompositumerkennung
- Kompositumzerlegung
- Mehrwortgruppenerkennung
- Wortrelationierung

Information Retrieval

Lingo

Mehrwortgruppenerkennung

- algorithmisch
- lexikalisch

Open source

...

alle Zeichenketten

unkontrollierte Zeichenketten

Volltextindexierung

-

(open

Abstract:

algorithmische

an

Anwendungsgrenzen

auf.

Automatische

automatischen

automatischen

basierten

Bei

Beispielen

Beitrag

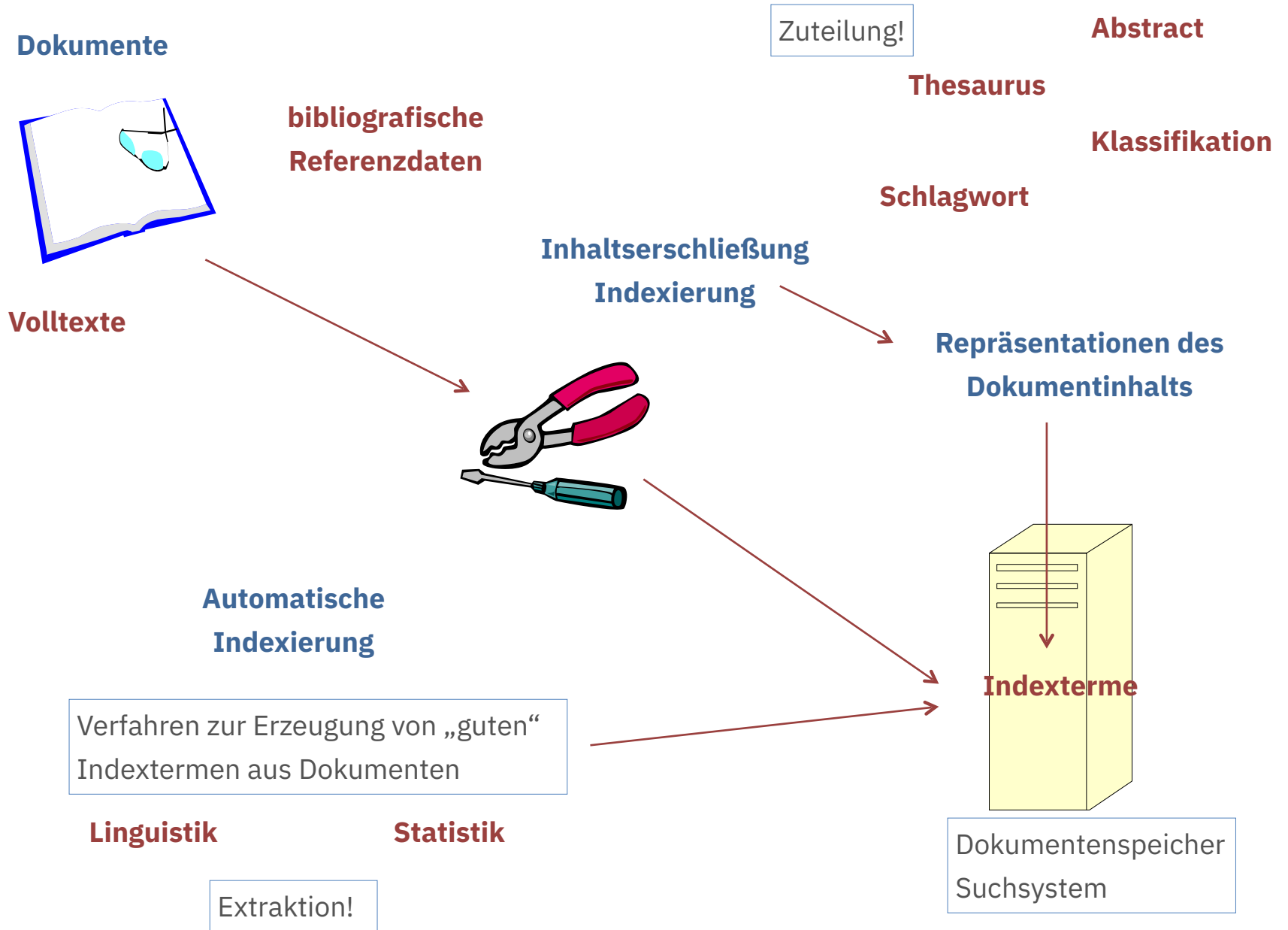
benannt.

...

Abb. 5.5 Gegenüberstellung manueller und automatischer Indexerstellung

Automatisches Indexieren

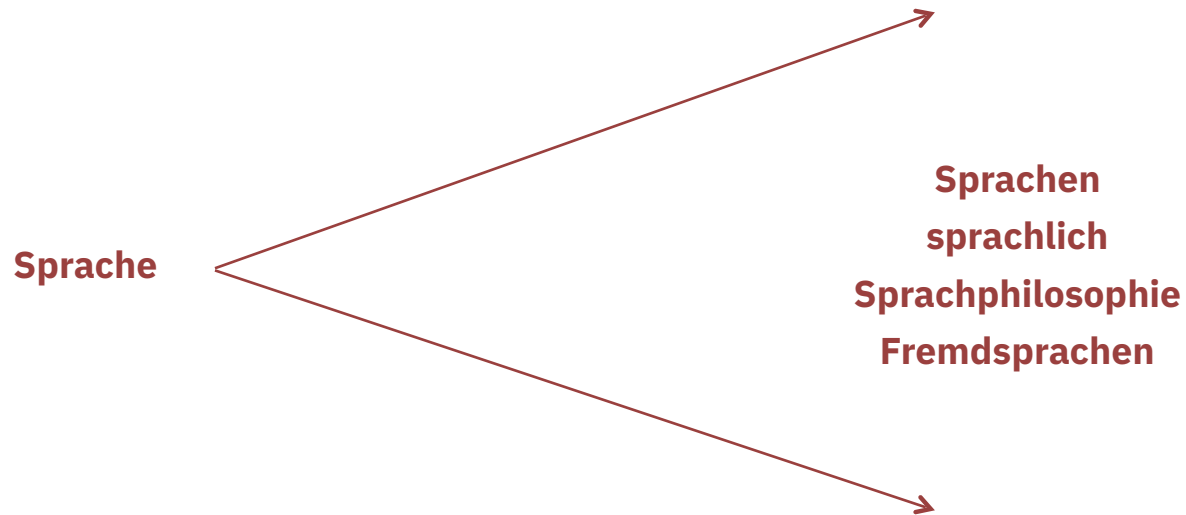
Indexieren und Automatisches Indexieren



Hypothese

„Wörter in Dokumenten sind selten gute Indexterme,
weil sie sprachlich zu verschieden sind.“

Verschiedenartigkeit von Dokument- und Suchsprache



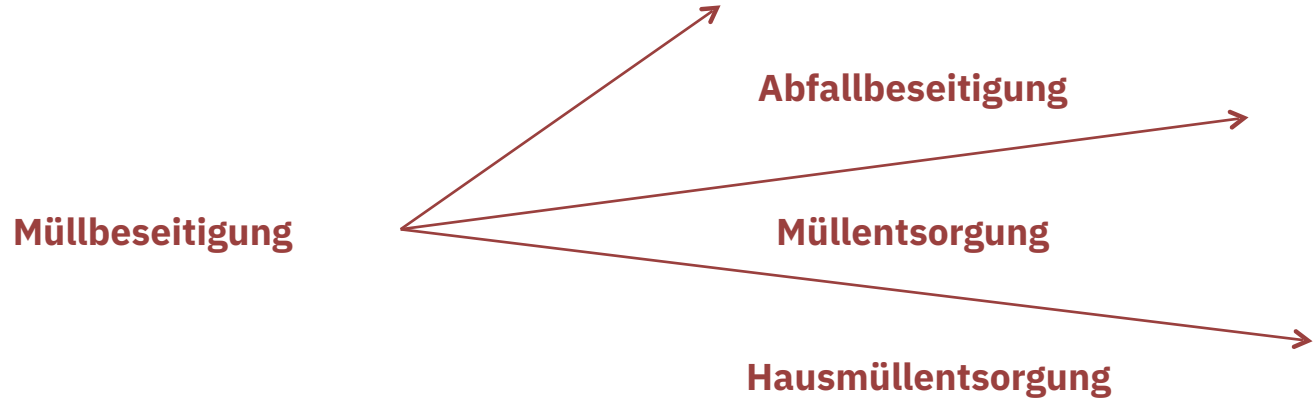
Lösung: Einsatz einer morphologischen Komponente

Wörter im Dokument

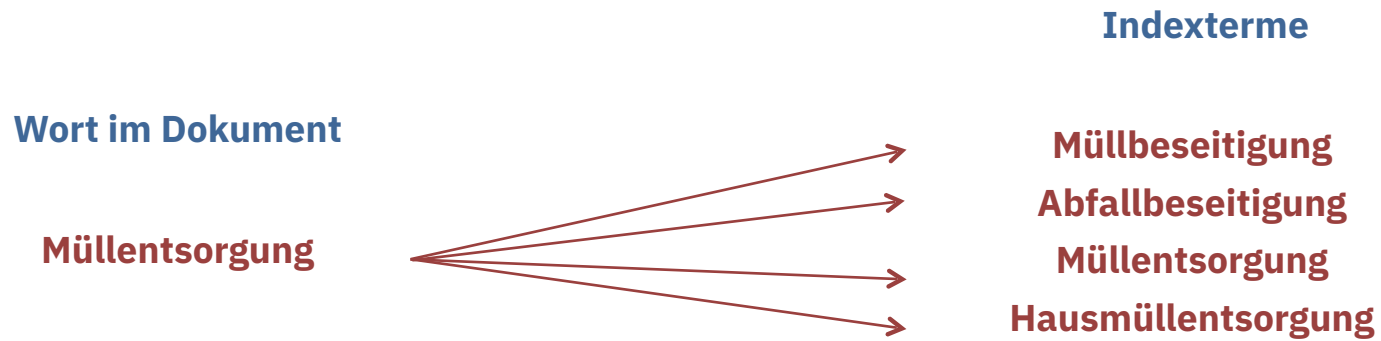
Indexterme



Problem: Suche im semantischen Umfeld



Lösung: Einbindung von Synonym- ggf. hierarchischen Relationen



Hypothese

„Nur Wörter mit bestimmten statistischen Merkmalen (Häufigkeitsmerkmalen) sind „gute“ Indexterme.“

Verfahren: Ermittlung von Worthäufigkeiten (z.B. TF, IDF)

Einsatzmöglichkeiten und Nutzen

Termgewichtung

**gewichtete
Indexterme**

Relevance Ranking

**Automatische
Klassifizierung**

**Notationen
eines
Klassifikationssystems**

**Systematische
Suche**

Dokumentclustering

**Mengen ähnlicher
Dokumente**

**Navigation in
großen Dokumentkollektionen**

Automatische Indexierung vs. Intellektuelle Erschließung

Ziele	<ul style="list-style-type: none"> • Retrievalverbesserung durch geeignete Indexterme • Recallerhöhung auf Stichwortbasis 	<ul style="list-style-type: none"> • einheitliche thematische Suche durch thematische, inhaltliche Zusammenführung • inhaltliche Repräsentation
Aufwand	<ul style="list-style-type: none"> • gering 	<ul style="list-style-type: none"> • hoch
Kosten	<ul style="list-style-type: none"> • gering 	<ul style="list-style-type: none"> • hoch
Ergebnisse	<ul style="list-style-type: none"> • Indexterme • vorhersagbar • beliebig reproduzierbar • durch wiederholte Indexierungen leicht zu verbessern 	<ul style="list-style-type: none"> • Schlagwörter, Deskriptoren, Notationen • menschliche Fehlerquote • endgültig • inhaltlich verlässlich
Qualität	<ul style="list-style-type: none"> • abhängig von Qualität der Stichwörter im Quelltext 	<ul style="list-style-type: none"> • hoch, bei konsistenter Erschließung
Nutzen	<ul style="list-style-type: none"> • für alle textbasierten Quelldokumente (auch erschlossene!) hoch • insb. für nicht erschlossene und nicht zu erschließende Dokumente hoch 	<ul style="list-style-type: none"> • bezogen auf die erschlossenen Dokumente hoch • in Kollektionen mit niedriger Erschließungsquote gering

Was tun?

- Linguistische Verfahren verbessern die Suchbarkeit allein auf der Basis von Stichwörtern. Aus nicht geeigneten bzw. nicht vorhandenen Stichwörtern lässt sich kein brauchbares Indexat machen!
- Statistische Verfahren bewerten die vorhandene Stichwortbasis – s.o.
- Klassische Erschließungsziele wie Zusammenführung von Gleichem, vollständiger Nachweis, zuverlässige und einheitliche inhaltliche Suche sind nur durch Erschließung, nicht durch automatische Indexierung zu erreichen.
- Eine Entscheidung für automatische Indexierung muss das Wissen um deren Leistungsfähigkeit insb. deren Grenzen der Leistungsfähigkeit berücksichtigen.
- Einer Entscheidung für automatische Indexierung sollte daher eine Zielbestimmung hinsichtlich der gewünschten Retrievalmöglichkeiten jetzt und in 10 Jahren vorausgehen. Entscheidungen gegen eine Erschließung sind nur mit erheblichen Konsistenzverlusten umkehrbar.

**Allgemein sind die Grenzen automatischer Erschließungsverfahren dort erreicht,
wo die Intelligenz beginnt.**

Automatische Schlagwortvergabe

Die „Automatische Schlagwortvergabe“ von Midos erzeugt auf der Basis von sog. Synonymlisten Deskriptoren zu Datensätzen. Dabei werden die Zeichenketten in den für das Verfahren ausgewählten Kategorien mit den Einträgen der Synonymliste verglichen und die erkannten bzw. erzeugten Deskriptoren in eine eigene Kategorie des Datensatzes geschrieben.

Thesaurus mit Allgemeindeskriptoren (deskr.mth)

MIDOSThesaurus-Viewer - DESKR

Suche nach

Hierarchische Liste | Alphabetische Listen | Beide | Relationen für Term

Liste

- Alle Terme
- Gleichordnendes Indexieren
- Grafische Methoden (BF)
- Güte des Suchergebnisses (BF)
- Güterklassifikation (BF)
- Güterklassifikationen (BF)
- Hardware
- Hardware (ENG)
- Heterogenitätsbehandlung**
- Heterogenitätsproblem (BF)
- Hierarchical relationship (ENG)
- Hierarchie
- Hierarchien (BF)
- Hierarchische Beziehung
- Hierarchische Beziehungen (BF)
- Hierarchy (ENG)
- Hilfsmittel des Information Retrieval (BF)
- Hilfstafl (BF)
- Historische Darstellung (BF)
- Homonym
- Homonyme (BF)
- Homonym (ENG)
- Human computer interface (ENG)
- Hyperlink
- Hyperlinks (BF)
- Hyperlink (ENG)

DE **Heterogenitätsbehandlung**

D Methoden, aus unterschiedlichen Erschließungsdaten einen gemeinsamen Suchprozess zu gestalten.

- TT [Katalog](#)
- BF [Heterogenitätsproblem](#)
- OB [Kataloganreicherung](#)
- VB [Citation Pearl Growing](#)
- VB [Semantische Interoperabilität](#)
- VB [Standardisierung](#)

Export

„deskriptor;synonym1;synonym2 ...“

Terme 1009 | Deskriptoren 281 | Verweise 2380 | Status OK | Datei C:\MIDOS-DB\LITERATUR\DESKR.mth

Wortliste (auto-sw.txt)

C:\gln-daten\wortlisten\auto-sw.txt

File Edit View Options

auto-sw.txt

Indexierung;Indexieren
 Indexierung
 Indexierungsverfahren;Indexierungsmethode
 Indexierungsprinzip
 Indexierungsprinzip;Indexierungsprinzipien
 Dokumentationssprache;Indexierungssprache
 Dokumentationssprache;Indexierungssprachen
 Retrievalstudie;Indexierungsstudie
 Retrievalstudie;Indexierungsstudien
 Indexierungsverfahren
 Indikatives Abstract
 Information
 Information Gateway
 Information Gateway;Information Gateways
 Information Management System
 Information Management System;Information Management Systeme
 Information Retrieval
 Modelle des Information Retrieval;Information Retrieval Modell
 Modelle des Information Retrieval;Information Retrieval Modelle
 Informationelle Selbstbestimmung;Informationelle Autonomie
 Informationskompetenz;Informationelle Kompetenz
 Informationelle Selbstbestimmung
 Informationsüberflutung;Informationelle Überflutung
 Informationsüberflutung;Informationeller Überfluss
 Information;Informationen
 Externalisierung;Informations-Externalisierung
 Information Management System;Informations-Ressourcen-System
 Informationsangebot
 Informationelle Selbstbestimmung;Informationsautonomie
 Informationsbedarf
 Informationsmittel;Informationsdatenbank

Dialog „Automatische Schlagwortvergabe“

The screenshot shows a dialog box titled 'Automatische Schlagwortvergabe' with the subtitle 'Auf der Basis von Positivlisten werden Texten Schlagwörter zugeordnet'. The dialog contains several input fields and buttons. Red arrows point from external text labels to specific elements in the dialog:

- zu bearbeitende Datenbank** points to the 'Quell-Datei' field containing 'C:\LIT-LINGO\MIDOS-LITERATUR\LITERATUR.DBM'.
- zu durchsuchende Felder** points to the 'Suchfelder' field containing 'TIT;ZUSTIT;ABS'.
- zu verwendende Positivliste** points to the 'Positivliste 1' field containing 'C:\lit-lingo\midos-literatur\synonym.WTX'.
- verschlagwortete Datenbank** points to the 'Ziel-Datei' field containing 'C:\LIT-LINGO\MIDOS-LITERATUR\literatur-01.dbm'.
- Feld für die Ergebnisse** points to the 'Zielfeld' field containing 'AUTOSW'.
- Blick in die Positivliste** points to the 'Anzeigen' button next to the 'Positivliste 1' field.
- Einstellungen für zusätzliche sprachliche Verarbeitung** points to the bottom section containing options for word length changes and language processing.

The dialog also features a 'Starten' button at the bottom left, a 'Ziel => Quelle' and 'Zieldatei löschen' button in the center, and a 'Schließen' button at the bottom right. A 'mehr' link is visible in the top right corner.

Testen!

Automatische Indexierung

Linguistisch basierte Verfahren

Indexliste

ixtrieve.fh-koeln.de/mwr/litie/mindexneu.pl 150%

Index zum Feld **Basic Index** (168957 Einträge) [alternative Liste](#)

Gehe zu: Suche: Ähnlich: suche alle

711	Suche
2	Suchecki
2	Sucheffizienz
1	Suchei
4	Sucheingabe
1	sucheinschränkende
1	Sucheinschränkung
10	Sucheinstieg
16	Sucheinstiege
3	Sucheinstiegen
1	Sucheintrag
1	Sucheinträge
1	Suchelement
357	Suchen
1	suchen-und-finden
20	Suchende
20	Suchenden
4	Suchender
1	Suchengine
7	Suchens
4	Sucher
1	Sucherei
1	Sucherergebnisse
2	Sucherfolg
1	Suchergänzungen
1	Suchergebnisse
24	Suchergebnis
1	Suchergebnis-Anordnung
1	Suchergebnis-Listen
1	Suchergebnis-Positionen
1	Suchergebnisfilter
5	Suchergebnisliste
1	Suchergebnislisten
116	Suchergebnisse
2	Suchergebnisseiten
41	Suchergebnissen
7	Suchergebnisses
5	Sucherlebnis
1	sucherweiternde
1	Sucherweiterung
12	Suchet

Verknüpfen mit oder (OR) und (AND) nicht (NOT)

Übernehmen Schließen

Stichwort - Wortform - Grundform

Zielvorstellung
 sprachlich begründete Reduzierung von
 Wortformen in Dokumenten (Stichwörter) auf
 Grundformen

Wortform

Suchen

Grundform

Suche

Wortform

Suchergebnisse

Grundform

Suchergebnis

Suche

Ergebnis

Wortform

sucheinschränkende

Grundform

sucheinschränkend

Sucheinschränkung

Suche

Einschränkung

Bestandteile der Sprache I

Phonem: kleinstes bedeutungsunterscheidendes Lautmerkmal

Morphem: kleinste bedeutungstragende Einheiten einer Sprache

Wortform: Erscheinungsformen von Wörtern in der Sprache; Zuordnung zur lexikalischen Einheit, z.B.:

Wort: bedeutungstragende Einheiten der Sprache, bestehend aus einzelnen Morphemen oder einer Kombination mehrerer Morpheme;

Wörter können im Satz ausgetauscht werden und Satzglieder bilden

Syntagma: mehrere Wörter können zu Wortgruppen, Wortfolgen oder Wortverbindungen verknüpft werden

Satz: Wortfolge mit mindestens einem Objekt (Subjekt) und einem Prädikat

Maus - Haus

Mantel - Hantel

Be-haus-ung, haus-en, Haus-ierer

**Haus, Häuser, Hauses, Häusern etc.
hausen, hausend**

Haus [Substantiv; Gebäude]

Maus [Substantiv; 1. Tier, 2. PC-Bediengerät]

hausen [Verb, umgspr. für wohnen]

Der Mond ist aus grünem Käse.

Der Mond ist aus gelbem Käse.

der gelbe Mond

juristische Person

Regeln für den Schlagwortkatalog

Hans schläft.

Hans schläft gerne in meiner Vorlesung.

Morphologie - Wörter und ihre Bestandteile

einfache Wörter (Simplizia)

Uhr (Kernmorphem)

Uhr-en (Kernmorphem und Flexionsmorphem)

Ableitungen (Derivationen)

Ver-bind-ung-en

**(Kernmorphem, Flexionsmorphem,
Wortbildungsmorphem)**

Komposita (mind. 2 Kernmorpheme)

Uhr-en-ver-gleich-s-test

Wortbildung erfolgt zum Beispiel durch

hinzufügen von Präfixen zum Wortstamm

ver-walt-en

P K F

hinzufügen von Suffixen zum Wortstamm

Ver-walt-ung

P K F

Ver-un-rein-ig-ung

P P K S S

Grundformerzeugung - Lemmatisierung

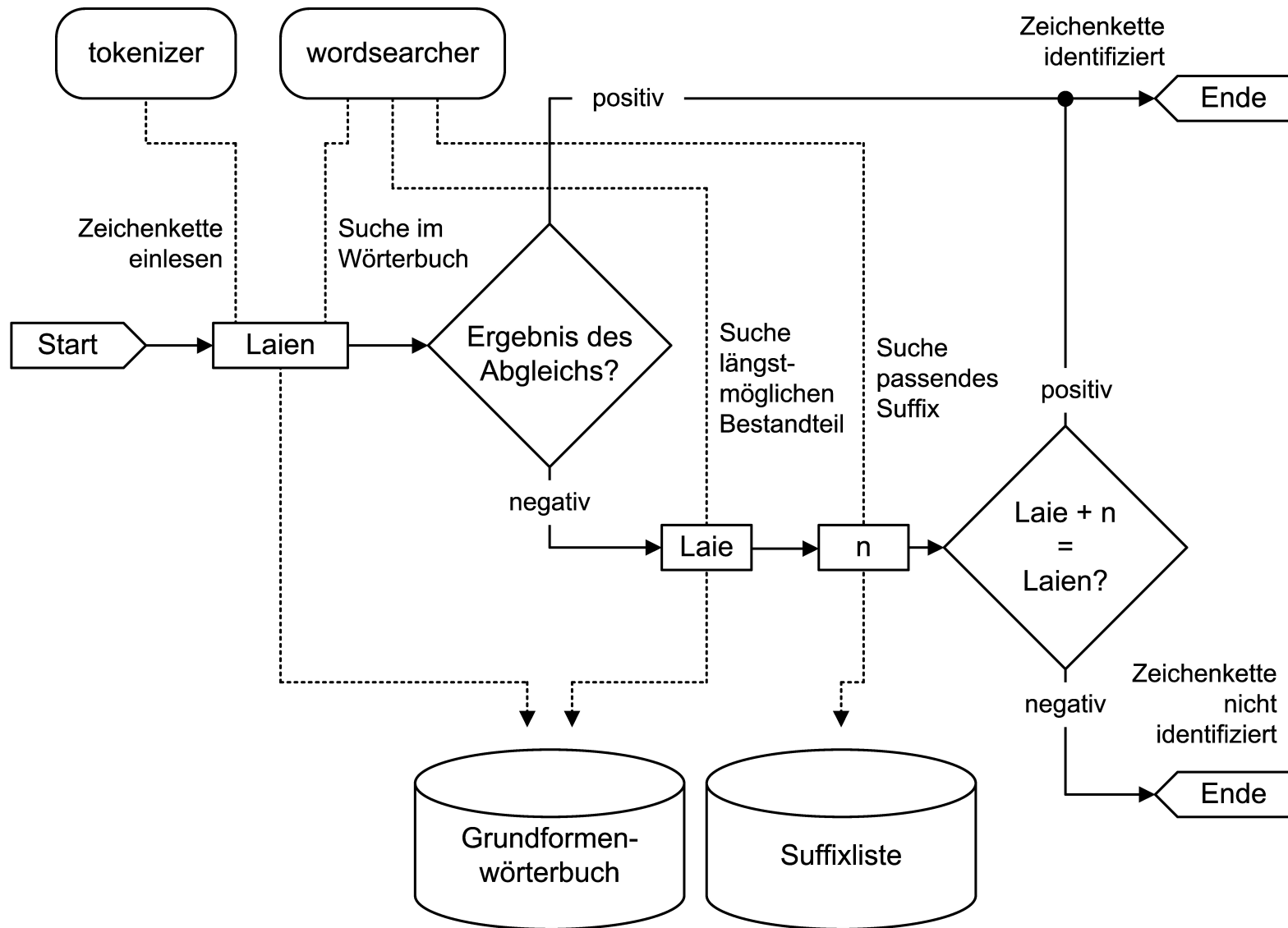


Abb. 5.14 Identifizierungsstrategie von Lingo

Tokenizer

Der tokenizer identifiziert Zeichenketten für den nachfolgenden Identifizierungsprozess (wordsearcher); Zeichen, die nicht Bestandteile von Wörtern sind, wie z.B. "&, \$, =, @", wirken als Begrenzer von Zeichenketten.

```
:Der/WORD/0/0:
:tokenizer/WORD/1/4:
:identifiziert/WORD/2/14:
:Zeichenketten/WORD/3/28:
:für/WORD/4/42:
:den/WORD/5/47:
:nachfolgenden/WORD/6/51:
:Identifizierungsprozess/WORD/7/65:
:wordsearcher/WORD/9/90:
(/OTHR/8/89:
:Zeichen/WORD/12/105:
)/OTHR/10/102:
;/PUNC/11/103:
:die/WORD/14/114:
:nicht/WORD/15/118:
,/PUNC/13/112:
:Bestandteile/WORD/16/124:
:von/WORD/17/137:
:Wörtern/WORD/18/141:
...
:sind/WORD/19/150:
:wie/WORD/21/156:
,/PUNC/20/154:
z.B/ABRV/22/160:
./PUNC/23/163:
"/OTHR/24/165:
:&/OTHR/25/166:
,/PUNC/26/167:
$/OTHR/27/169:
,/PUNC/28/170:
=/OTHR/29/172:
,/PUNC/30/173:
:wirken/WORD/34/179:
:@/OTHR/31/175:
"/OTHR/32/176:
:als/WORD/35/186:
,/PUNC/33/177:
:Begrenzer/WORD/36/190:
:von/WORD/37/200:
./PUNC/39/217:
```

```
# tokenizer rules:
# SPAC = \s+
# NUMS = [+ -]?(?:\d{4,}|\d{1,3}(?:\.\d{3,3})*) (?:\.| (?:,\d+)??)
# URLS = (?:www\.|mailto:|(?:news|https?|ftps?)://|\S+?[\._]\S+?@\S+?\.)\S+
# ABRV = (?: (?: (?:CHAR)+\.)+ ) (?:CHAR)+
# WORD = ALNUM(?:-*ALNUM)*
# PUNC = [! , . : ; ? ; ! ; ]
# OTHR = [ - "$ % & ' ( ) * + \ - / < = > @ \ [ \ \ \ ] ^ _ { | } ~ ¢ £ ¤ ¥ ¦ § ¨ © « ¬ ® ¯ ° ± ² ³ ´ µ ¶ · ¸ ¹ º » ¼ ½ ¾ × ÷ ]
# HELP = \S*
```

Wordsearcher

Wörterbücher von *Lingo* sind reine Textdateien

Wortform

```
lahmend=lahmend #a lahmen #v
lähmend=lähmend #a lähmen #v
lahmen=lahm #a lahmen #v
lähmen=lähmen #v
lahmgelegt=lahmgelegt #a lahmlegen #v
lahm=lahm #a lahmen #v
lahmlegend=lahmlegend #a lahmlegen #v
lähmung=lähmung #s.f
lahmzulegen=lahmzulegen #v
lahnung=lahnung #s.f
lahr=lahr #e.f
laib=laib #s.m
laichen=laichen #v
laie=laie #s.m
laienhaft=laienhaft #a
laientum=laientum #s.n
laizistisch=laizistisch #a
lakai=lakai #s.m
lake=lake #s.f
laken=laken #s.n
lakonisch=lakonisch #a
lakritze=lakritze #s.f
laktase=laktase #s.f lact
laktation=laktation #s.f
```

Grundform(en)

Wortklasse

Genus

```
# Systemwörterbuch lingo-dic.txt
# enthält deutsche Grundformen mit Wortklassen
#
# #s = Substantiv
# #a = Adjektiv
# #v = Verbe
# #e = Eigename
# #t = Takeitasis
# #w = Wortform
```

Identifizierungsstrategie

Token

Laien

Abgleich mit dem Wörterbuch

lahr=lahr #e.f
laib=laib #s.m
laichen=laichen #v
laie=laie #s.m
laienhaft=laienhaft #a
laientum=laientum #s.n
laizistisch=laizistisch #a
lakai=lakai #s.m
lake=lake #s.f
laken=laken #s.n
lakonisch=lakonisch #a
lakritze=lakritze #s.f
laktase=laktase #s.f lactase #s.f
laktation=laktation #s.f

längster
übereinstimmender
Eintrag

laie

Überprüfung:
falls „laien“ eine
zulässige Wortform der
Grundform „laie“ ist,
muss „n“ erlaubtes
Suffix für die Wortklasse
„s“ sein

Bestätigung

Ergebnis

<Laien = [(laie/s/m)]>

suffix:

Suffixliste, Stand: 30-06-2005

Suffixklasse: s = Substantiv, a = Adjektiv, v = Verb, e = Eigenwort, f = Fugung

Suffixe je Klasse: <suffix>['/'<ersetzung>][<suffix>['/'<ersetzung>]]"

- [s, 'e en er ern es **n** s se sen ses']

- [a, 'este ste ster sten stes ester estes esten e em en er ere eren erer eres es erem']

- [v, 'e/en en/en est/en et/en st/en t/en te/en ten/en eten/en ete/en etest/en s']

- [e, 's']

- [f, 's n e en es er ch/che ch/chen']

Kompositumerkennung

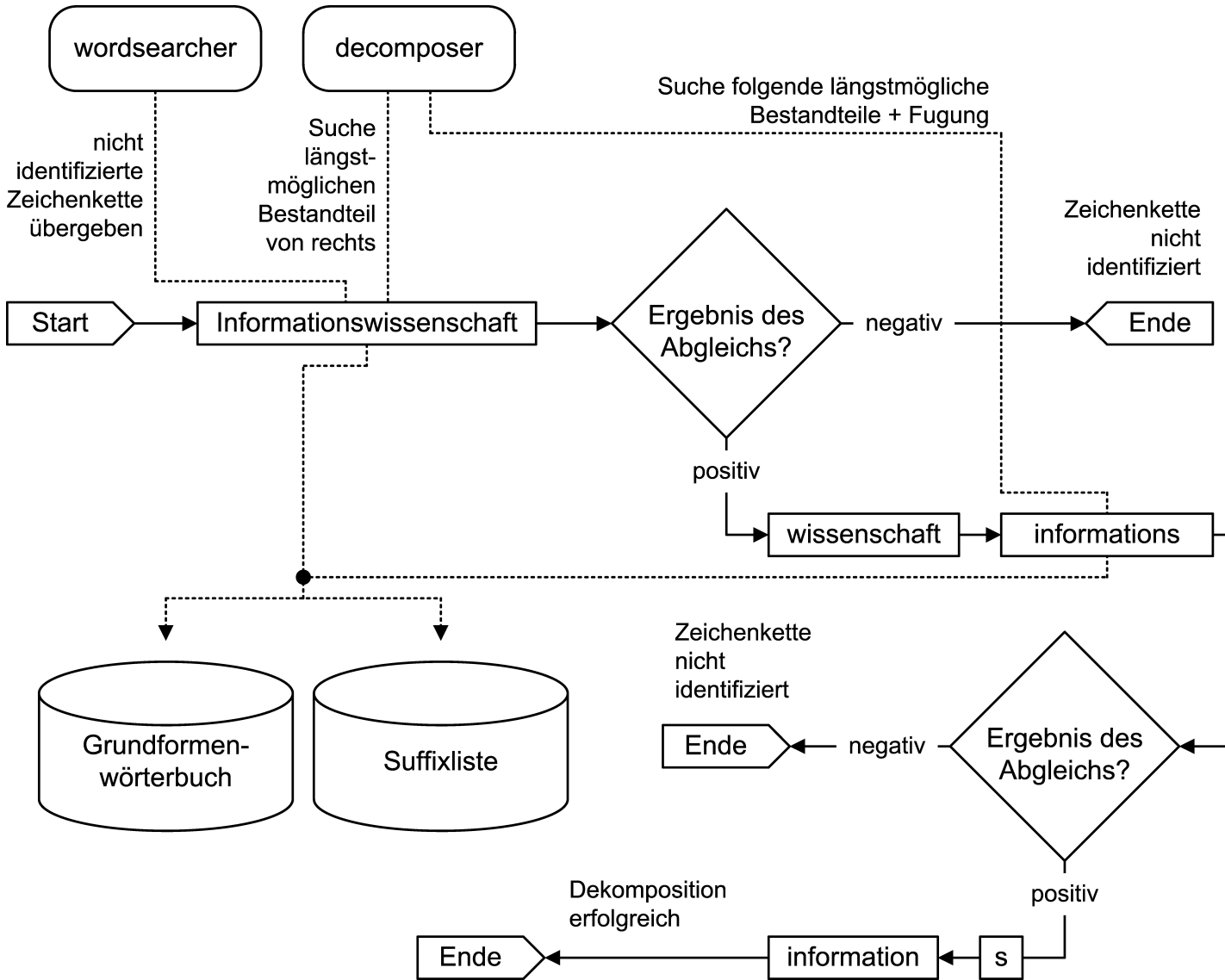


Abb. 5.16 Identifizierung von Komposita mit *Lingo*

Token

Informationswirtschaft

Informationswirtschaft

Suche nach längst-möglichem Bestandteil von rechts

Informations

positiv - Suche nach vorderem Teil

Information-s

Bestätigung

Ergebnis

Überprüfung auf zulässige Endung im Kompositum (Fugung)

<informationswirtschaft|COM = [(informationswirtschaft/k), (information/s+/f), (wirtschaft/s+/f)]>

Abgleich mit dem Wörterbuch

nicht im Wörterbuch

- informant=informant #s.f
- informatiker=informatiker #s.m
- informatik=informatik #s.f
- information=information #s.f**
- informationsverarbeitend=informationsverarbeitend #a
- informativ=informativ #a**
- informativisch=informativisch #a

- wirtschaftlichkeit=wirtschaftlichkeit #s.f
- wirtschaftlich=wirtschaftlich #a
- wirtschaft=wirtschaft #s.f**
- wirt=wirt #s.m
- wir=wir #w

suffix:

- # Suffixliste, Stand: 30-06-2005
- # Suffixklasse: s = Substantiv, a = Adjektiv, v = Verb, e = Eigenwort, f = Fugung
- # Suffixe je Klasse: "<suffix>['/'<ersetzung>][<suffix>['/'<ersetzung>]]"
- [s, 'e en er ern es n s se sen ses']
- [a, 'este ste ster sten stes ester estes esten e em en er ere eren erer eres es erem']
- [v, 'e/en en/en est/en et/en st/en t/en te/en ten/en eten/en ete/en etest/en s']
- [↓, 's']
- [**f**, **s** n e en es er ch/che ch/chen']

Konfiguration (lit-lingo-lem.cfg)

Analyseteil

```
#  
# Lingo-Konfiguration  
# Version "lem" - nur tokenizer, word_searcher, decomposer, vectorfilter  
# 2021-02-16 / le  
---  
meeting:  
  
  attendees:  
  
  #####  
  # Text bereitstellen  
  #  
  
  # Angegebene Datei zeilenweise einlesen und verarbeiten  
  - text_reader: { files: $(files), progress: true }  
  
  #####  
  # Inhalte verarbeiten  
  #  
  
  # Zeile in einzelnen Sinnbestandteile (Token) zerlegen  
  - tokenizer: { }  
  
  # Verbleibende Token im Wörterbuch suchen  
  - word_searcher: { source: 'usr-dic sys-dic' mode: first }  
  
  # Nicht erkannte Wörter auf Kompositum testen  
  - decomposer: { source: 'usr-dic sys-dic' out: res }
```

verwendete Wörterbücher

**Verarbeitung nach
erster Identifizierung
beenden**

Ausgabe des

Ergebnisses als „res“

Konfiguration (lit-lingo-lem.cfg)

Ausgabeteil

attendee für die Ausgabe
der „log“-Datei

Filtern der Ergebnisse in
„res“ und Sortierung

Ausgabe der Ergebnisse
in Datei

```
#####
```

```
# Ergebnisse ausgeben
```

```
#
```

Einlesen des Ergebnisses aus dem Analyseteil

```
# Erstelle Datei mit Endung .log für Datenstrom
```

```
- debug_filter: { in: res, prompt: 'lex:) ' }
```

```
- text_writer: { ext: log, sep: "\n" }
```

```
# Erstelle Datei mit Endung .non für nicht erkannte Wörter
```

```
- vector_filter { in: res, lexicals: '\?' sort: term_abs }
```

```
- text_writer: { ext: non, sep: "\n" }
```

Filtereinstellungen: Wortklassen

```
# Erstelle Datei mit Endung .vec für erkannte Indexterme
```

```
- vector_filter: { in: res, lexicals: '^[ksave]$\'} }
```

```
- text_writer: { ext: vec, sep: "\n" }
```

```
# Erstelle Datei mit Endung .ven für erkannte Indexterme mit absoluter Häufigkeit
```

```
- vector_filter: { in: res, lexicals: '^[ksavem]$', sort: term_abs }
```

```
- text_writer: { ext: ven, sep: "\n" }
```

**Dateiendung jeder Eintrag
auf neuer Zeile**

**Sortierung der
ausgegebenen
Ergebnisse:
absolute
Termhäufigkeit**

Mehrwortererkennung

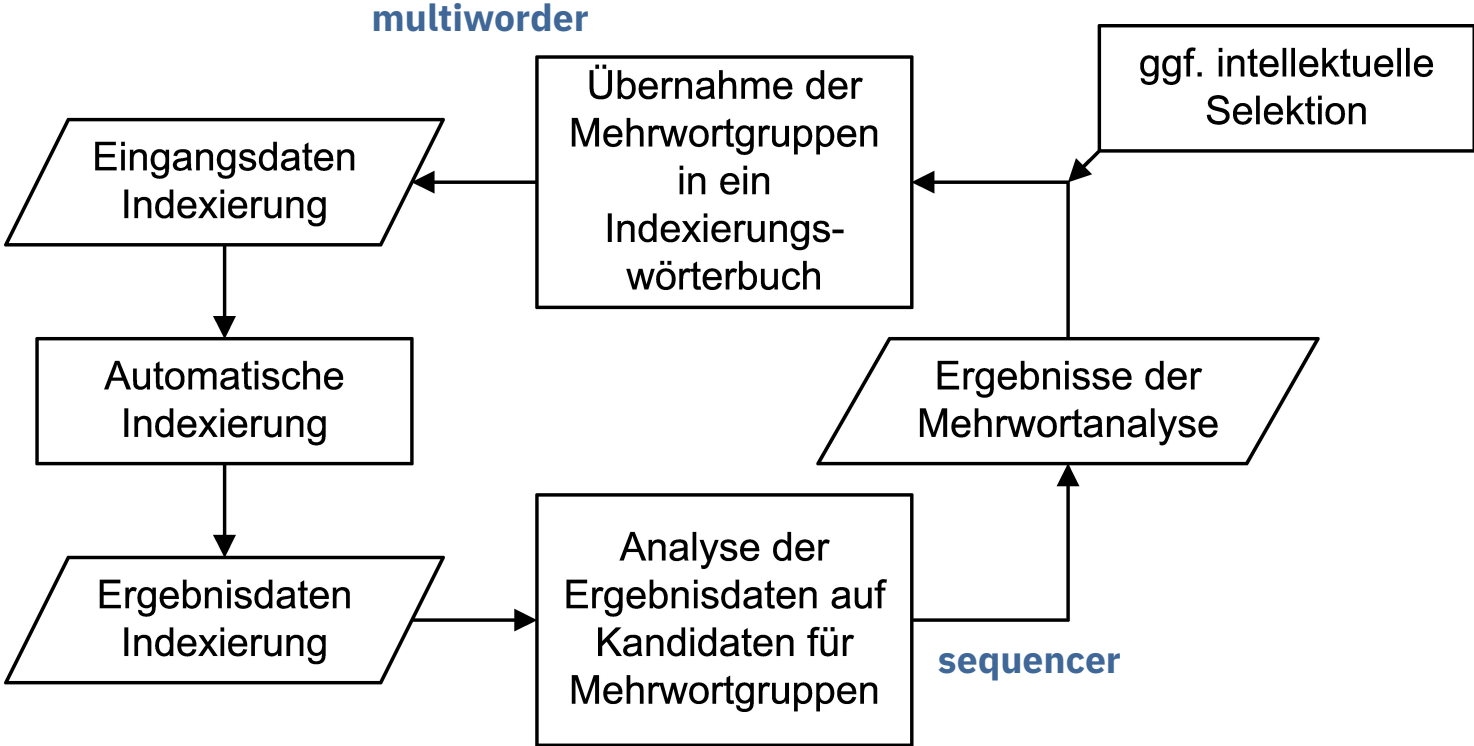


Abb. 5.17 Ablaufschema der Mehrworterkennung

lit-lingo-sem.cfg

```
---
meeting:

  attendees:

    #####
    # Text bereitstellen
    #

    # Angegebene Datei zeilenweise einlesen und verarbeiten
    - text_reader:      { files: $(files), progress: true }

    #####
    # Inhalte verarbeiten
    #

    # Zeile in einzelnen Sinnbestandteile (Token) zerlegen
    - tokenizer:        { }

    # Verbleibende Token im Wörterbuch suchen
    - word_searcher:   { source: 'usr-dic sys-dic', mode: first }

    # Nicht erkannte Wörter auf Kompositum testen
    - decomposer:      { source: 'usr-dic sys-dic' }

    # Mehrwortgruppen im Strom erkennen
    - multi_worder:    { source: sys-mul }

    # Wortsequenzen anhand von Regeln identifizieren
    - sequencer:       { stopper: 'PUNC,OTHR', out: res }
```

wörterbuchbasierte Mehrworterkennung

regelbasierte Mehrworterkennung

Sequencer - algorithmische Identifizierung von Mehrwortgruppen

Konfigurationsdatei „de.lang“

```
sequencer:
    sequences: [ [AS, '2, 1'], [AK, '2, 1'], [AAK, '3, 1 2'], [AAS, '3, 1 2'] ]
```

	Wortklassenmuster	Ausgabe
automatisch akquisition -- akquisition, automatisch		
automatisch bestimmung -- bestimmung, automatisch		
automatisch dokumenterschließung -- dokumenterschließung, automatisch		AK
automatisch dokumentklassifikation -- dokumentklassifikation, automatisch		
automatisch erschließung -- erschließung, automatisch		AS
automatisch gruppierung -- gruppierung, automatisch		
automatisch indexieren -- indexieren, automatisch		
automatisch indexierung -- indexierung, automatisch		
automatisch indexierungssystem -- indexierungssystem, automatisch		
automatisch inhaltlich erschließung -- erschließung, automatisch inhaltlich		AAS
automatisch inhaltserschließung -- inhaltserschließung, automatisch		
automatisch klassifikation -- klassifikation, automatisch		
automatisch maschine -- maschine, automatisch		
automatisch methode -- methode, automatisch		
automatisch recherche -- recherche, automatisch		
automatisch selektion -- selektion, automatisch		
automatisch semantisch klassifikation -- klassifikation, automatisch semantisch		
automatisch thematisch textklassifikation -- textklassifikation, automatisch thematisch		AAK
automatisch verfahren -- verfahren, automatisch		
automatisch vollindexierung -- vollindexierung, automatisch		
automatisch wortformenreduktion -- wortformenreduktion, automatisch		
automatisch übersetzungssystem -- übersetzungssystem, automatisch		

Multiworder - wörterbuchbasierte Identifizierung von Mehrwortgruppen

lingo-mul.txt

```
automatische bildverarbeitung
automatische datenverarbeitung
automatische differentiation
automatische dokumentation
automatische fahrtzielführung
automatische fahrzeugführung
automatische flugsteuerung
automatische folge
automatische handlungsplanung
automatische inhaltsanalyse
automatische justierung
automatische lastabhängige bremskraftregelung
```

Lingo-Wörterbücher sind zwar reine Textdateien, werden aber vor der Indexierung zu **Indexdateien** (*sdbm*) konvertiert, um die Verarbeitung zu beschleunigen. Dabei kann eine Vorverarbeitung stattfinden, z.B. eine Lemmatisierung der Einträge, um auch Varianten erkennen zu können.

automatisch lastabhängig bremskraftregelung #/m

erkennt:

```
automatische lastabhängige bremskraftregelung
automatische lastabhängige bremskraftregelungen
automatischen lastabhängigen bremskraftregelung
automatischen lastabhängigen bremskraftregelungen
```

--- de.lang

```
language:
  name: 'Deutsch'
```

dictionary: **interner Name**

Vorverarbeitung: Lemmatisierung

databases:

Quelldatei

Format der Einträge

Wortklasse

```
# System dictionaries
```

```
sys-dic: { name: de/lingo-dic.txt, txt-format: WordClass, separator: '=' }
```

```
sys-abk: { name: de/lingo-abk.txt, txt-format: WordClass, separator: '=' }
```

```
sys-syn: { name: de/lingo-syn.txt, txt-format: KeyValue, separator: '=', def-wc: y }
```

```
sys-mul: { name: de/lingo-mul.txt, txt-format: SingleWord, use-lex: sys-dic, def-wc: m }
```

Synonymer - Zuteilung von Synonymen als Indexterme

	gnd-syn-lingo.txt
VIR*Visual Information Retrieval	
Viröse Vergilbung*Vergilbungskrankheit	
Viröse Zweignekrose*Viröse Zweignekrose	
Viraginität*Virilismus	
Viral Marketing*Viral Marketing	
Virale hämorrhagische Septikämie*Hämorrhagische Virusseptikämie	
Virale Infektion*Virusinfektion	
Virales Insektizid*Virales Insektizid	
Virales Marketing*Viral Marketing	
virales onkogen*virales onkogen	
Viramid*Ribavirin	
Virashaivas*Wiraschaiwa	
Virasoro-Algebra*Virasoro-Algebra	
Virazol*Ribavirin	
Virelai*Virelai	
Viren / Hülle*Virushülle	
Viren / Inaktivierung*Virusinaktivierung	
Viren*Viren	
Virenhemmendes Mittel*Virostatikum	
Virenneutralisation*Virusneutralisation	
virensscanner*virensscanner	
Virenschutzpaket DOS*VSPdos	
Virenschutzpaket Windows 3.0*VSPwin	
Virenschutzsoftware*Virenschutzsoftware	

„Viraginität“ im Text erzeugt zusätzlich „Virilismus“

Synonyme können Mehrwortgruppen sein, daher zunächst Mehrworterkennung, danach Synonymgenerierung

„Hilfseintrag“ zur Identifizierung von Deskriptoren

Eigennamen

GND-Besonderheiten: Kombinationen

mehrwortige Eigennamen

```
# User dictionaries de.lang Synonym*Vorzugsbenennung Wortklasse  
usr-dic: { name: de/user-dic.txt, txt-format: WordClass, separator: '=' }  
usr-syn: { name: de/gnd-syn-lingo.txt, txt-format: KeyValue, separator: '*', def-wc: y }
```

Alternatives Wörterbuchformat: MultiValue

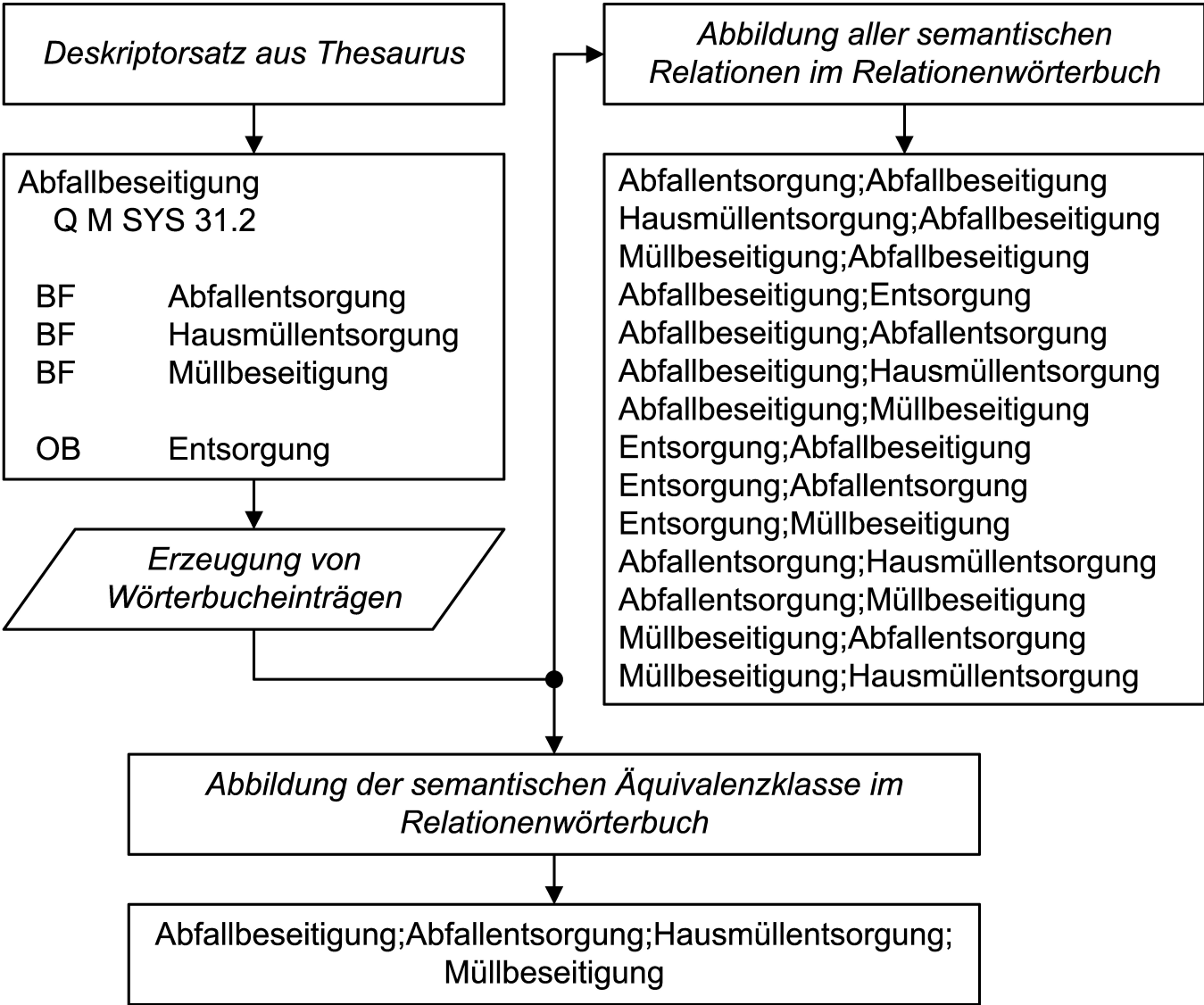


Abb. 5.19 Aufbau der Relationenwörterbücher für Synonyme

Indexieren einer Dokumentenkollektion

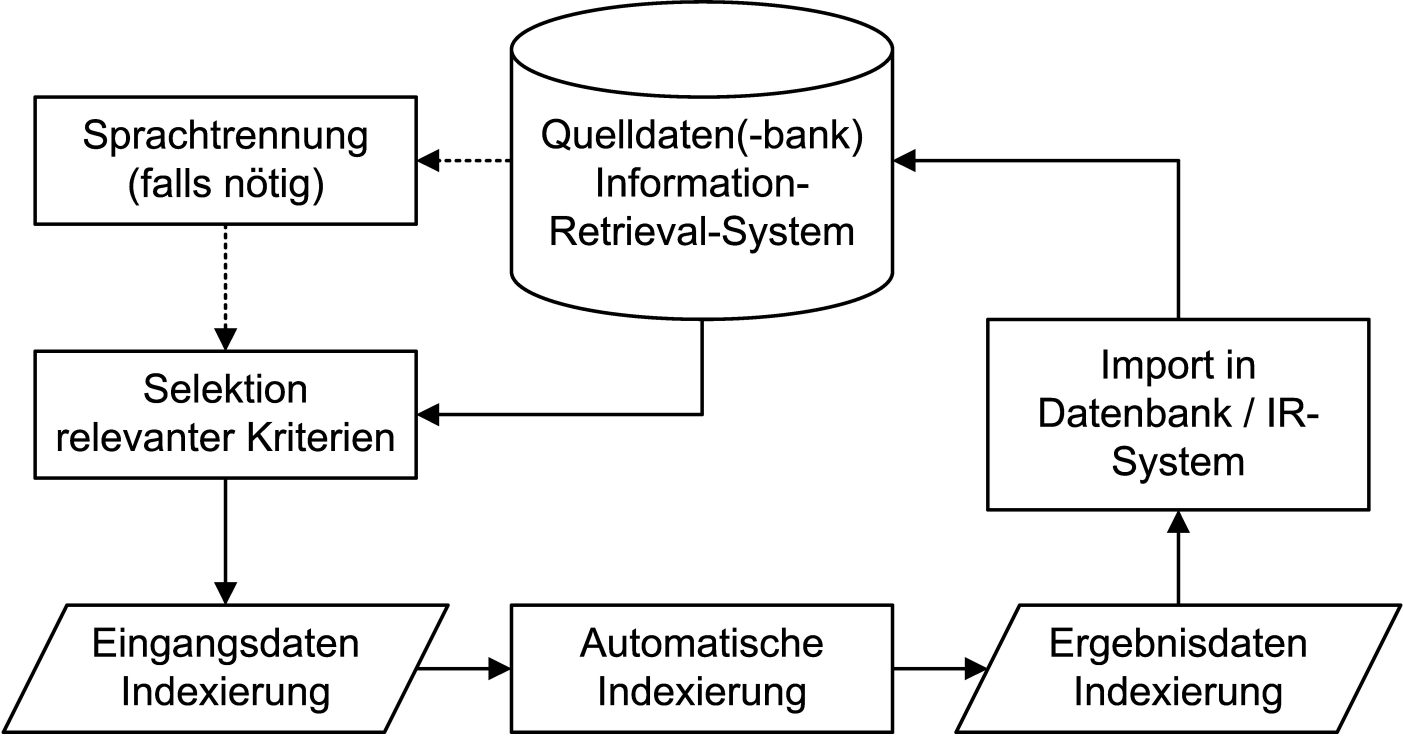


Abb. 5.20 Ablaufschema eines Indexierungslaufs

Export der Daten

MIDOS 6 - Vollanzeige
Anzeige Anzeige erweitern Ausgabe

literatur.dbm 60 / 100 / 100

Ioannidis, G.; Hermes, T.; Miene, A.: Wie kommt das Bild in die Datenbank?
Inhaltsbasierte Analyse von Bildern und Videos.
In: Information - Wissenschaft und Praxis. 53(2002) H.1, S.15-21.

Abstract: Die verfügbare multimediale Information nimmt stetig zu, nicht zuletzt durch die Tag für Tag wachsende Zahl an neuer Information im Internet. Damit man dieser Flut Herr werden und diese Information wieder abrufbar machen kann, muss sie annotiert und geeignet in Datenbanken abgelegt werden. Hier besteht das Problem der manuellen Annotation, das einerseits durch die Ermüdung aufgrund der Routinearbeit und andererseits durch die Subjektivität des Annotierenden zu Fehlern in der Annotation führen kann. Unterstützende Systeme, die dem Dokumentar genau diese Routinearbeit abnehmen, können hier bis zu einem gewissen Grad Abhilfe schaffen. Die wissenschaftliche Erschließung von beispielsweise filmbeiträgen wird der Dokumentar zwar immer noch selbst machen müssen und auch sollen, aber die Erkennung und Dokumentation von sog. Einstellungsgrenzen kann durchaus automatisch mit Unterstützung eines Rechners geschehen. In diesem Beitrag zeigen wir anhand von Projekten, die wir durchgeführt haben, wie weit diese Unterstützung des Dokumentars bei der Annotation von Bildern und Videos gehen kann.

Formate
neu
ändern
S-Form
0 Vollanzeige
1 Tabelle
2 Montiertes Aus
3 lingo-export
4 lingo-export-tt

Datensatz im montierten Ausgabeformat

MIDOS 6 - Vollanzeige
Anzeige Anzeige erweitern Ausgabe

literatur.dbm 60 / 100 / 100

[100245.]
Wie kommt das Bild in die Datenbank?
Inhaltsbasierte Analyse von Bildern und Videos
Die verfügbare multimediale Information nimmt stetig zu, nicht zuletzt durch die Tag für Tag wachsende Zahl an neuer Information im Internet. Damit man dieser Flut Herr werden und diese Information wieder abrufbar machen kann, muss sie annotiert und geeignet in Datenbanken abgelegt werden. Hier besteht das Problem der manuellen Annotation, das einerseits durch die Ermüdung aufgrund der Routinearbeit und andererseits durch die Subjektivität des Annotierenden zu Fehlern in der Annotation führen kann. Unterstützende Systeme, die dem Dokumentar genau diese Routinearbeit abnehmen, können hier bis zu einem gewissen Grad Abhilfe schaffen. Die wissenschaftliche Erschließung von beispielsweise filmbeiträgen wird der Dokumentar zwar immer noch selbst machen müssen und auch sollen, aber die Erkennung und Dokumentation von sog. Einstellungsgrenzen kann durchaus automatisch mit Unterstützung eines Rechners geschehen. In diesem Beitrag zeigen wir anhand von Projekten, die wir durchgeführt haben, wie weit diese Unterstützung des Dokumentars bei der Annotation von Bildern und Videos gehen kann.

Formate
neu
ändern
S-Form
0 Vollanzeige
1 Tabelle
2 Montiertes Aus
3 lingo-export
4 lingo-export-tt

Datensatz im Format „lingo-export“

Export-Dialog in Midos

Treffer speichern (Export)

von Treffer 1 bisTreffer 100 Satzabstand 1

Datei wählen

Dateiname export.txt

Speicherformate

- HTML
- RTF
- Text (CSV)
- Text (Wahlfelder)
- Text (Ausgabeformate)
- MIDOS
- MIDOS (Wahlfelder)

Text (Ausgabeformat)

- nur aktuellen Treffer speichern
- an existierender Datei anhängen
- keine Abfrage bei Dateixistenz
- Fenster schließen nach Export
- UTF8

Speichern Ansicht Schließen

```
#  
# Lingo-Konfiguration für Dateien im LIR-Format  
#  
# Gebräuchliche Patterns sind  
#  
# "\021(\d+\-\d+)\022"  
# "\^[(\d+)\.]"  
#  
---
```

„LIR-Modus“: Indexieren einer Datei, die aus einzelnen Datensätzen besteht

Erkennungsmuster für die Identnummer eines Datensatzes

meeting:

attendees:

#####

Text bereitstellen

#

„records: true“ - datensatzweise Indexierung

Angegebene Datei zeilenweise einlesen und verarbeiten

- text_reader: { files: \$(files), records: true, progress: true }

[100019.] exportierte Datensätze

Halbautomatische Volltextanalyse, Datenbankaufbau und Document Retrieval

[100022.]

Nutzung von Klassifikationssystemen zur verbesserten Beschreibung, Organisation und Suche von Internetressourcen

100019*datenbankaufbau|halbautomatisch|retrieval|volltextanalyse

Indexierungsergebnis

100022*beschreibung|beschreibung, verbessert|internetressource|klassifikationssystem|nutzung|organisation|suche|suchen|verbessert

Identnummer

100019*datenbankaufbau|halbautomatisch|retrieval|volltextanalyse

Indexierungsergebnis „export.vec“

100022*beschreibung|beschreibung, verbessert internetressource|klassifikationssystem|nutzung|organisation|suche|suchen|verbessert

Indexterm

Trenner: |

Trenner für Felder ID, Lingo

Job-Datei für Import in Midos

Job-Kommandofolge für MWUPDATE (Achtung !!! Datei nur mit MWUPDATE bearbeiten)

UTF8TOANSI

C:\lit-lingo\lingo-work\txt\export.vec

C:\lit-lingo\lingo-work\txt\lingo.csv

1. Umbenennen in „csv“ und konvertieren in UTF-8

UTF8TOANSI

C:\lit-lingo\lingo-work\txt\export.syn

C:\lit-lingo\lingo-work\txt\lingo-gnd.csv

READAUSTAUSCH

C:\lit-lingo\lingo-work\txt\lingo.csv

C:\lit-lingo\lingo-work\txt\lingo.dbm

FELD:DOKNR;LINGO TZ:* SZ:^ nohead noEZ

READAUSTAUSCH

C:\lit-lingo\lingo-work\txt\lingo-gnd.csv

C:\lit-lingo\lingo-work\txt\lingo-gnd.dbm

FELD:DOKNR;LINGOGND TZ:* SZ:^ nohead noEZ

MISCHDATNEU

C:\lit-lingo\midos-literatur\literatur-02.dbm

C:\lit-lingo\lingo-work\txt\lingo.dbm

C:\lit-lingo\midos-literatur\literatur-03.dbm KEYFELD:DOKNR

MISCHDATNEU

C:\lit-lingo\midos-literatur\literatur-03.dbm

C:\lit-lingo\lingo-work\txt\lingo-gnd.dbm

C:\lit-lingo\midos-literatur\literatur-04.dbm KEYFELD:DOKNR

2. Konvertieren von „csv“ in Midos-Datenbankformat „dbm“

2. Einmischen der Indexierungsergebnisse im „dbm“-Format in die Datensätze der Datenbank und umbenennen der Datenbank

Ergebnis des Midos-Jobs

MIDOS 6 - Vollanzeige

Anzeige Anzeige erweitern Ausgabe

literatur-04.dbm 18 / 100 / 100

0: Vollanzeige

Dokumentnummer 100033
Verfasser Oeltjen, W.
Titel Dokumentenstrukturen manipulieren und visualisieren
Zusatz HST über das Arbeiten mit der logischen Struktur
Quelle Herausforderungen an die Wissensorganisation: Visualisierung, multimediale Dokumente, Internetstrukturen. 5. Tagung der Deutschen Sektion der Internationalen Gesellschaft für Wissensorganisation Berlin, 07.-10. Oktober 1997. Hrsg.: H. Czap u.a.
Verlagsort Würzburg
Verlag Ergon Verlag
Erscheinungsjahr 1998
Umfang S.151-156
Serie Fortschritte in der Wissensorganisation; Bd.5
Sprache deutsch
Dokumenttyp in
Abstract Thema dieses Beitrages sind Dokumentenstrukturen und zwar aus zwei Blickrichtungen: aus der Sicht der Autoren, die ein Dokument mit Computerunterstützung erstellen und die Dokumentenstruktur manipulieren und aus der Sicht der Lesenden, die ein Dokument lesen und die Struktur des Dokumentes wahrnehmen. Bei der Dokumentenstruktur wird unterschieden zwischen der logischen Struktur und der grafischen Struktur eines Dokumentes. Diese Trennung ermöglicht das Manipulieren und Visualisieren der logischen Struktur. Welche Bedeutung das für die Autoren und für die Benutzenden des Dokumentes hat, soll in dem Beitrag u.a. am Beispiel der Auszeichnungssprache HTML, der Dokumentenbeschreibungssprache des World-Wide Web, erörtert werden.
Deskriptoren Deskriptive Auszeichnungssprache; Semantisches Datenmodell
Objekte SGML; HTML
Automatische SW-Zuteilung Visualisierung
Automatische GND-Zuteilung Dokument; Struktur; Manipulation; und; Visualisierung

Lingo-Indexterme arbeiten; auszeichnungssprache; autor; bedeutung; beispiel; beitrag; benutzend; computerunterstützung; dokument; dokumentenbeschreibungssprache; dokumentenstruktur; einen; ermöglichen; ermöglicht; erstellen; erörtern; erörtert; grafisch; graphisch; haben; lese; lesen; lesend; logisch; manipulieren; sicht; sichten; soll; sollen; struktur; struktur, grafisch; struktur, logisch; trennung; unterscheiden; visualisieren; wahrnehmen; werden

Lingo-GND arbeiten; aus; auszeichnungssprache; autor; bedeutung; beispiel; beitrag; computerunterstütztes verfahren; dokument; ein; haben; in; intelligentes netz; lesen; mit; sicht; sichten; sollen; struktur; trennung; und; wahrnehmen; werden; zwischen

eingemischte Indexterme

Formate
neu
ändern
S-Form
0 Vollanzeige
1 Tabelle
2 Montiertes Aus
3 lingo-export
4 lingo-export-tit

Stemming-Verfahren - algorithmische Lösungen für die sprachliche Vereinheitlichung

Voraussetzung

Die Angleichung der Wörter an die Grammatik erfolgt regelhaft mit einem überschaubaren Vorrat an Regeln und einer nicht zu großen Zahl an Ausnahmen. Typischer Einsatzbereich ist die Verarbeitung der englischen Sprache.

famil-ies → **famil-y**

satisf-ied → **satisfy**

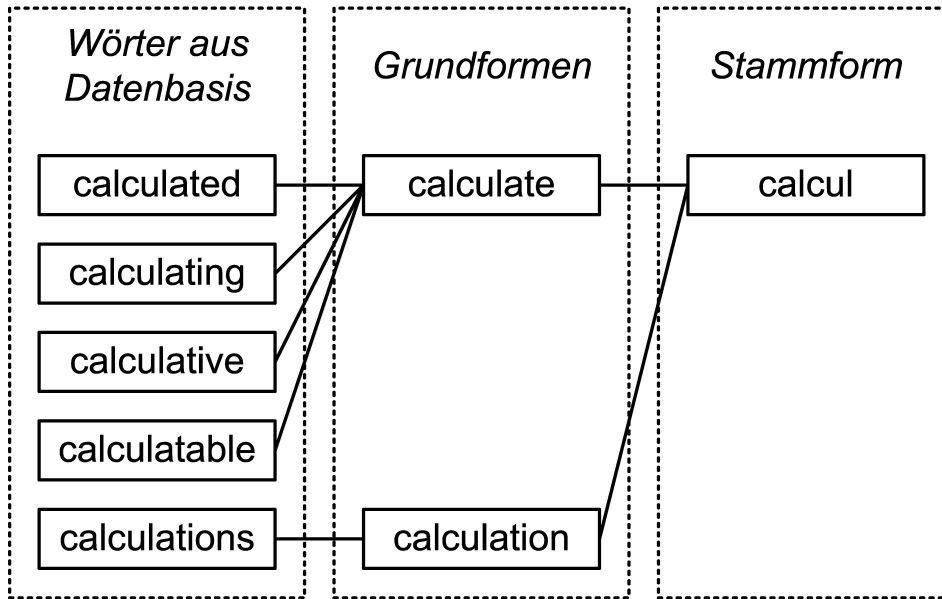
index-ing → **index**

retriev-ed → **retrieve**

1.	IES	-->	Y	
2.	ES	-->	_	[wenn *O / CH / SH / SS / ZZ / X vorausgehen]
3.	S	-->	_	[wenn * / E / %Y / %O / OA / EA vorausgehen]
4.	IES'	-->	Y	
	ES'	-->	_	
	S'	-->	_	
5.	'S	-->	_	
	'	-->	_	
6.	ING	-->	_	[wenn ** / % / X vorausgehen]
	ING	-->	E	[wenn %* vorausgehen]
7.	IED	-->	Y	
8.	ED	-->	_	[wenn ** / % / X vorausgehen]
	ED	-->	E	[wenn %* vorausgehen]

% = alle Vokale und Y
* = alle Konsonanten
_ = Tilgung
/ = ODER

Abb. 5.26 Kuhlen-Algorithmus



Grundformen sind vollständige Wörter,
Stammformen nicht

Grundformen sind taugliche Indexterme,
Stammformen nicht – falls Stämme Indexterm
werden, muss auch der Suchterm „gestemmt“
werden.

Abb. 5.27 Stammformreduktion für das Englische

Overstemming

- winner → win-ner
- winning → win-ning
- wine → win-e

Understemming

- divide → divid-e
- dividing → divid-ing
- division → divis-ion

Porter: „An algorithm for suffix stripping“

These may all be represented by the single form

`[C]VCVC ... [V]` (V=Vokal und C=Konsonant)

where the square brackets denote arbitrary presence of their contents.

Using `(VC){m}` to denote VC repeated m times, this may again be written as

`[C](VC){m}[V]`. **m = Zahl der Wiederholungen**

m will be called the \measure\ of any word or word part when represented in this form. The case m = 0 covers the null word. Here are some examples:

- m=0 TR, EE, TREE, Y, BY. **Beispiele**
- m=1 TROUBLE, OATS, TREES, IVY.
- m=2 TROUBLES, PRIVATE, OATEN, ORRERY.

The \rules\ for removing a suffix will be given in the form

`(condition) S1 -> S2` **Regeln für Suffixentfernung**

This means that if a word ends with the suffix S1, and the stem before S1 satisfies the given condition, S1 is replaced by S2. The condition is usually given in terms of m, e.g.

`(m > 1) EMENT ->`

zusätzliche Bedingungen

The 'condition' part may also contain the following:

- *S - the stem ends with S (and similarly for the other letters).
- *v* - the stem contains a vowel.
- *d - the stem ends with a double consonant (e.g. -TT, -SS).
- *o - the stem ends cvc, where the second c is not W, X or Y (e.g. -WIL, -HOP).

Abfolge im Algorithmus

Step 1a

SSES -> SS	caresses -> caress
S ->	cats -> cat

Step 1b

(m>0) EED -> EE	agreed -> agree
(*v*) ED ->	plastered -> plaster
(*v*) ING ->	motoring -> motor

If the second or third of the rules in Step 1b is successful, the following is done:

AT -> ATE	conflat(ed) -> conflate
BL -> BLE	troubl(ed) -> trouble
IZ -> IZE	siz(ed) -> size

Step 2

`(m>0) IZATION -> IZE` vietnamization -> vietnamize

Step 3

`(m>0) ALIZE -> AL` formalize -> formal

Step 4

`(m>1) AL ->` revival -> reviv

generalizations

- generalizations -> generalization # S ->
- generalization -> generalize # (m>0) IZATION -> IZE
- generalize -> general # (m>0) ALIZE -> AL
- general -> gener # (m>1) AL

gener

Porter, Martin F.: „An algorithm for suffix stripping“, in: Program 14/3 (1980), S. 130–137.

<https://tartarus.org/martin/PorterStemmer/def.txt>

Automatische Indexierung

Statistisch basierte Verfahren

Allgemeine Überlegungen

- Es besteht ein Zusammenhang zwischen der Auftretenshäufigkeit von Wörtern und deren Bedeutung für das Retrieval.
- Wichtig sind solche Wörter, die
 - Dokumente hinreichend signifikant vertreten und gleichzeitig
 - von nicht-relevanten Dokumenten trennen.

Zipf's law

Worthäufigkeit * **Häufigkeitsrang** = Konstante

Worthäufigkeit: Auftretenshäufigkeit eines Wortes in einer Kollektion

Häufigkeitsrang: Position in der nach Häufigkeit sortierten Liste

Beispiel:

1.	Häufigstes Wort	10.000
2.	Zweithäufigstes Wort	5.000
3.	Dritthäufigstes Wort	3.300
...		
10.000.	Zehntausend...	1

Vermutungen

- hochfrequente Wörter sind schlechte Suchbegriffe, weil sie unspezifisch sind
- niedrigfrequente Wörter sind schlechte Suchbegriffe, weil sie wahrscheinlich nicht zum Vokabular des Nutzers gehören und/oder autorenspezifisch sind

Wortverteilung in den Kollektionen von TREC-1 (1993)

Quelle	WSJ	AP	ZIFF	FR	DOE
Größe des Korpus in MB	295	266	251	258	190
Mittelwert Wörter/Datensatz	182	353	181	313	82
Verschiedene Wörter im Korpus	156.000	198.000	174.000	126.000	186.000
Dokumenthäufigkeit = 1	65.000	90.000	86.000	59.000	96.000
Mittelwert bei Dokumenthäuf. > 1	199	174	165	106	159

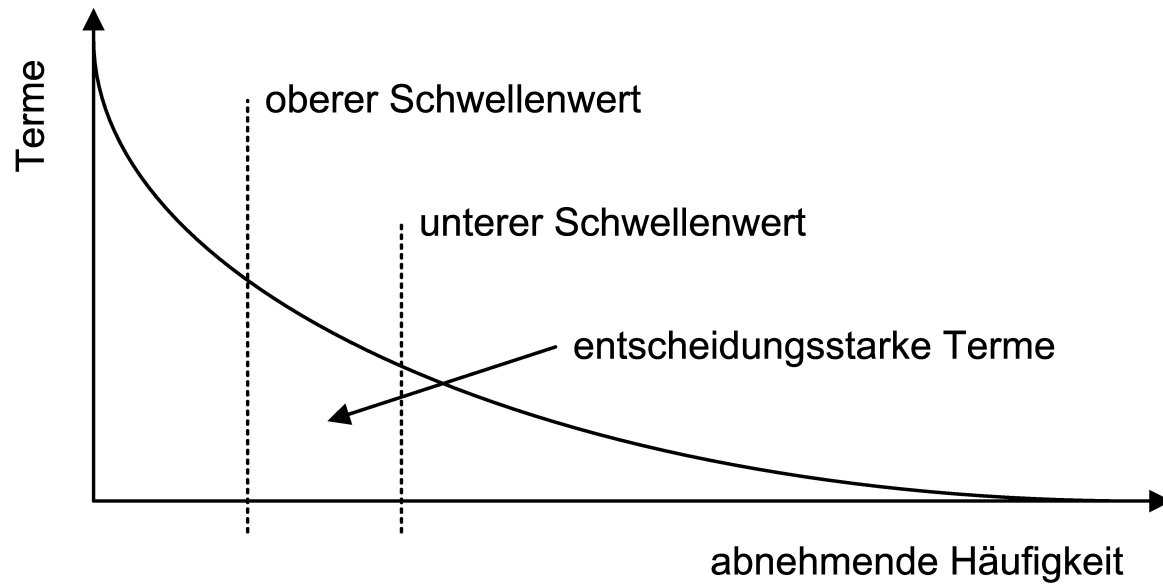


Abb. 5.6 Verteilung von Termhäufigkeiten

Termgewichtung

(1) Einfache Termhäufigkeit (TF)

„Je häufiger ein Term in einem Dokument vorkommt, umso wichtiger ist er für dieses Dokument.“

Termhäufigkeit = Häufigkeit Term je Dokument

(2) Relative Termhäufigkeit (WDF)

„Die einfache Termhäufigkeit bevorzugt lange Dokumente, interessant ist also die Termhäufigkeit in Relation zur Länge des Dokuments.“

WDF = Häufigkeit Term je Dokument / Gesamtzahl Terme im Dokument

(3) Inverse Dokumenthäufigkeit (IDF)

„Terme, die in wenigen Dokumenten vorkommen, sind wichtiger als Terme, die in sehr vielen oder fast allen Dokumenten vorkommen (Bsp. Funktionswörter)“

IDF = 1 / Häufigkeit der Dokumente je Term

(4) Termhäufigkeit * Inverse Dokumenthäufigkeit (TF*IDF)

„Terme, die in relativ wenigen Dokumenten relativ häufig vorkommen, sind wichtig.“

TF*IDF = Termhäufigkeit bzw. WDF / Dokumenthäufigkeit

LIR – "Lehr- und Lernsystem Information Retrieval"

New query | Query

Help | LIR home

Searching database Literatur zur Informationerschießung Change

Query[?]:

Options[?]: Regular expression End-truncated Case-sensitive

Ranking[?]: Kascade komplex Threshold: @100

Custom[?]:

Search

Reset

Literatur zur Informationerschießung (42.740 records)

Search tips:

- Query will be performed as "ranked query"[?]
- Try regular expressions[?]
- +/-TERM = TERM *must / must not* occur
- TERM^N = TERM's weight gets multiplied by NUMBER (=> *boost factor*)
- Custom ranking[?]:
 tf = Term frequency
 df = Document frequency
 N = Number of documents in the collection
- Threshold:
NUMBER(%) = Display only documents with at least a weight of NUMBER (percent of the maximum weight)
 @NUMBER(%) = Display top NUMBER (percent) documents

Sample query: [+indexierung_ranking.*^1.5 automatisch^-1.5 -thesaurus](#)

New query | Query

Help | LIR home

LIR – "Lehr- und Lernsystem Information Retrieval"

[New query](#) | [Query](#) | [Results](#)

! [Help](#) | [LIR home](#)

Searching database "Literatur zur Informationserschließung"

Query:

Options: Regular expression End-truncated Case-sensitive

Ranking: Threshold:

Art der Termgewichtung

Displaying results 1 - 100 (100 of [241](#) results)

Select ranking: and threshold:

Relevance Feedback

1. (231.) [00059](#) @10.2192 (4.6068) > Ranking-Experimente mit gewichteter Indexierung [**indexierung**]
2. (217.) [00282](#) @10.2192 (4.898) > Syntax und Gewichtung in Informationssprachen : Ein Fortschrittsbericht über präzisere Indexierung und Computer-Suche [**indexierung**]
3. (136.) [00670](#) @10.2192 (5.1434) > Transparente Indexierungsstrukturen im Fach Literaturwissenschaft [**indexierung**]
4. (215.) [00833](#) @10.2192 (4.9658) > DIN 31623: Indexierung zur inhaltlichen Erschließung von Dokumenten : T.1: Begriffe, Grundlagen; T.2: Gleichordnende Indexierung mit Deskriptoren; T.3: Syntaktische Indexierung mit Deskriptoren [**indexierung**]
5. (232.) [00930](#) @10.2192 (4.6068) > Entwicklung und Fortschritt bei Klassifikation und Indexierung [**indexierung**]
6. (126.) [00954](#) @10.2192 (5.1458) > Automatische Indexierung zwischen Forschung und Anwendung [**indexierung**]
7. (160.) [01461](#) @10.2192 (5.1324) > Probabilistische Modelle in Information-Retrieval-Systemen [**indexierung**]
8. (44.) [01490](#) @10.2192 (5.1718) > Wäre es nicht langsam Zeit, die Informationstechnologie in der bibliothekarischen Sacherschließung etwas erster zu nehmen? : ein Wort zur RSWK [**indexierung**]
9. (127.) [01528](#) @10.2192 (5.1452) > –Die Pilotstudie "DB-Thesaurus" : Allgemeiner Thesaurus für Bibliotheken [**indexierung**]
10. (192.) [01533](#) @10.2192 (5.1092) > Verbindliche versus freie Indexierung [**indexierung**]
11. (218.) [01535](#) @10.2192 (4.898) > LyberWorld : eine 3D-basierte Benutzerschnittstelle für die computerunterstützte Informationssuche in Dokumentmengen [**indexierung**]
12. (158.) [01553](#) @10.2192 (5.135) > Zukunftsperspektiven der Klassifikation und Indexierung [**indexierung**]
13. (67.) [01555](#) @10.2192 (5.1646) > –Die Funktion von semantischen Kategorien in Indexierungssprachen und bei der Indexierung [**indexierung**]
14. (233.) [02006](#) @10.2192 (4.6068) > Grundriß eines Thesaurus als funktionsfähiges Hilfsmittel für Indexierung und Recherche [**indexierung**]
15. (151.) [02319](#) @10.2192 (5.1374) > Automatische Indexierung: Entwicklung und Perspektiven [**indexierung**]

LIR – "Lehr- und Lernsystem Information Retrieval"

New query | Query | Results

Help | LIR home

Searching database "Literatur zur Informationserschließung"

Query:

Options: Regular expression End-truncated Case-sensitive

Ranking: Threshold:

Neues Ranking mit verändertem Gewichtungsalgorithmus

Displaying results 1 - 100 (100 of 241 results)

Select ranking: and threshold:

1. (139.) [31099](#) @5.193 (10.2192) > Weltkongress Bibliothek und Information, 72. IFLA-Generalkonferenz in Seoul, Korea : Aus den Veranstaltungen der Division IV Bibliographic Control, der Core Activities ICABS und UNIMARC sowie der Information Technology Section [**indexierung**]
2. (155.) [32695](#) @5.1928 (10.2192) > Weltkongress Bibliothek und Information : 73. IFLA-Generalkonferenz in Durban, Südafrika. Aus den Veranstaltungen der Division IV Bibliographic Control, der Core Activities (CABS und UNIMARC sowie der Information Technology Section [**indexierung**]
3. (118.) [29305](#) @5.1918 (10.2192) > Weltkongress Bibliothek und Information, 71. IFLA-Generalkonferenz in Oslo, Norwegen : Aus den Veranstaltungen der Division IV Bibliographic Control, der Core Activities ICABS und UNIMARC sowie der Information Technology Section [**indexierung**]
4. (124.) [30144](#) @5.1912 (10.2192) > -Die Volltextabfrage und das Alleinstellungsmerkmal des physischen Buches [**indexierung**]
5. (143.) [31276](#) @5.1908 (10.2192) > Durchführung von Digitalisierungsprojekten in Bibliotheken [**indexierung**]
6. (108.) [27877](#) @5.1898 (10.2192) > Indexierung von Online-Katalogen : Ein gemeinsames Konzept der ALEPH-Anwender in Berlin [**indexierung**]
7. (162.) [33828](#) @5.1896 (10.2192) > Neue Plattform für Indexierung und Retrieval : Studenten erstellen die Webseite »iXtrieve« [**indexierung**]
8. (176.) [35150](#) @5.1896 (10.2192) > Zwischenbilanz Collaborative Catalog Enrichment [**indexierung**]
9. (138.) [31082](#) @5.188 (10.2192) > Automatische Indexierung des Reallexikons zur Deutschen Kunstgeschichte [**indexierung**]
10. (142.) [31218](#) @5.186 (10.2192) > Information und Sprache : Beiträge zu Informationswissenschaft, Computerlinguistik, Bibliothekswesen und verwandten Fächern. Festschrift für Harald H. Zimmermann. Herausgegeben von Ilse Harms, Heinz-Dirk Luckhardt und Hans W. Giessen [**indexierung**]
11. (99.) [27319](#) @5.1858 (10.2192) > Suchmaschine mit Mehrwert : Mirago [**indexierung**]
12. (136.) [31079](#) @5.1844 (10.2192) > Indexieren mit AUTINDEX [**indexierung**]
13. (95.) [27185](#) @5.1842 (10.2192) > Sacherschließung - wir müssen sie (uns) leisten! : Vorträge im Rahmen der 28. Jahrestagung der Gesellschaft für Klassifikation, Universität Dortmund 9. bis 11. Mai 2004 [**indexierung**]
14. (107.) [27876](#) @5.1832 (10.2192) > Weltkongress Bibliothek und Information, 70. IFLA-Generalkonferenz in Buenos Aires : Aus den Veranstaltungen der Division IV Bibliographic Control, der Core Activities ICABS und UNIMARC sowie der Information Technology Section [**indexierung**]
15. (119.) [29306](#) @5.1826 (10.2192) > Weltkongress Bibliothek und Information, 71. IFLA-Generalkonferenz in Oslo, Norwegen : Aus den Veranstaltungen der Division IV Bibliographic Control, der Core Activities ICABS und UNIMARC sowie der Information Technology Section [**indexierung**]

LIR – "Lehr- und Lernsystem Information Retrieval"

New query | Query | Results

Help | LIR home

Searching database "Literatur zur Informationserschließung"

Query:

Options: Regular expression End-truncated Case-sensitive

Ranking: Threshold:

Relevance Feedback über Markierung geeigneter Treffer und Ähnlichkeitssuche

Displaying results 1 - 100 (100 of 241 results)

Select ranking: and threshold:

1. (139.) [31099](#) @5.193 (10.2192) > Weltkongress Bibliothek und Information, 72. IFLA-Generalkonferenz in Seoul, Korea : Aus den Veranstaltungen der Division IV Bibliographic Control, der Core Activities ICABS und UNIMARC sowie der Information Technology Section [**indexierung**]
2. (155.) [32695](#) @5.1928 (10.2192) > Weltkongress Bibliothek und Information : 73. IFLA-Generalkonferenz in Durban, Südafrika. Aus den Veranstaltungen der Division IV Bibliographic Control, der Core Activities (CABS und UNIMARC sowie der Information Technology Section [**indexierung**]
3. (118.) [29305](#) @5.1918 (10.2192) > Weltkongress Bibliothek und Information, 71. IFLA-Generalkonferenz in Oslo, Norwegen : Aus den Veranstaltungen der Division IV Bibliographic Control, der Core Activities ICABS und UNIMARC sowie der Information Technology Section [**indexierung**]
4. (124.) [30144](#) @5.1912 (10.2192) > -Die Volltextabfrage und das Alleinstellungsmerkmal des physischen Buches [**indexierung**]
5. (143.) [31276](#) @5.1908 (10.2192) > Durchführung von Digitalisierungsprojekten in Bibliotheken [**indexierung**]
6. (108.) [27877](#) @5.1898 (10.2192) > Indexierung von Online-Katalogen : Ein gemeinsames Konzept der ALEPH-Anwender in Berlin [**indexierung**]
7. (162.) [33828](#) @5.1896 (10.2192) > Neue Plattform für Indexierung und Retrieval : Studenten erstellen die Webseite »iXtrieve« [**indexierung**]
8. (176.) [35150](#) @5.1896 (10.2192) > Zwischenbilanz Collaborative Catalog Enrichment [**indexierung**]

Relevance Feedback

Sukzessive Annäherung an die ideale Treffermenge durch

- Auswahl geeigneter Dokumente in der Treffermenge
- erneute Suche des Systems nach „ähnlichen“ Dokumenten
(bspw. durch eine Suche mit gemeinsamen hochgewichteten Indextermen)

LIR – "Lehr- und Lernsystem Information Retrieval"

New query | Query | Results | Records

Help | LIR home

Searching database "Literatur zur Informationserschließung" **neue Suche**

Query:

Options: Regular expression End-truncated Case-sensitive

Ranking: Threshold:

Relevance Feedback

Displaying results 1 - 100 (100 of 295 results)

Select ranking: and threshold:

neue Treffermenge durch Ähnlichkeitssuche

1. [35150](#) @55.1416 > Zwischenbilanz Collaborative Catalog Enrichment [**basis-technologie bibliotheksverbänden datenbank-management-system indexierungsergebnis inhaltsverzeichnis kataloganreicherung perspektivenwechsel retrieval-technologie sammlungsorganisation sprachverarbeitungskonzept**]
2. [31276](#) @55.0066 > Durchführung von Digitalisierungsprojekten in Bibliotheken [**bibliothekskatalog derivationsmorphologie digitalisierungsprojekt erschließungsmöglichkeit flexionsmorphologie indexierungssprache indexierungsverfahren indexterminologie rechtschreibprüfung unregelmäßigkeit**]
3. [30144](#) @54.2986 > ¬Die Volltextabfrage und das Alleinstellungsmerkmal des physischen Buches [**begriffsrecherche online-suchabfrage qualitätssteigerung registererstellung texterschließung trefferqualität verweisungs begriff volltextabfragen volltextretrieval wissenstheoretisch**]
4. [35151](#) @16.5346 > Vom Katalog zur Bibliothek : Zwischenschritt und Zwischenstand "Kataloganreicherung" [**basis-technologie datenbank-management-system kataloganreicherung**]
5. [36716](#) @16.4342 > Kataloganreicherung auch für Reihe B der Deutschen Nationalbibliografie [**bibliotheksverbänden inhaltsverzeichnis kataloganreicherung**]
6. [31041](#) @16.4222 > Kosten und Nutzung der Optimierung von Erschließung [**indexierungsverfahren inhaltsverzeichnis kataloganreicherung**]
7. [34220](#) @16.3558 > Kataloganreicherungsdaten der Deutschen Nationalbibliothek für Dritte zugänglich [**bibliothekskatalog inhaltsverzeichnis kataloganreicherung**]
8. [34488](#) @16.3556 > Kataloganreicherungsdaten der Deutschen Nationalbibliothek für Dritte zugänglich [**bibliothekskatalog inhaltsverzeichnis kataloganreicherung**]
9. [28558](#) @16.3538 > Aus der 48. Sitzung der Arbeitsgemeinschaft der Verbundsysteme am 12. und 13. November 2004 in Göttingen [**bibliothekskatalog inhaltsverzeichnis kataloganreicherung**]
10. [27434](#) @16.3496 > In Bibliothekskatalogen "googlen" : Integration von Inhaltsverzeichnissen, Volltexten und WEB-Ressourcen in Bibliothekskataloge [**bibliothekskatalog inhaltsverzeichnis kataloganreicherung**]

Help contents

Quick tips	– What are a few <i>quick tips</i> ?
Ranked query	– What's a <i>ranked query</i> ?
Algorithms	– Which ranking/weighting <i>algorithms</i> are supported?
Regular expressions	– What are <i>regular expressions</i> ?
Options	– What do these <i>options</i> mean?

Quick tips

- The query will be performed as "[ranked query](#)" (or "best match query", as opposed to "exact match query").
- There's no support for (hierarchical) grouping by parentheses; phrases ("...") are supported though.
- [Regular expressions](#) are supported!
- Terms beginning with a + sign **must** occur.
- Terms beginning with a – sign **must not** occur.
- You can modify ("boost") a term's weight by appending "^" plus value to the term which will be multiplied with the term's initial weight (negative values are possible as well).
- Sample query: `+indexierung_ranking.*^1.5 automatisch^-1.5 -thesaurus`
Will find documents containing "indexierung" and preferably "ranking" (the ".*" indicates end-truncation) but not "thesaurus", any occurrence of "automatisch" will decrease the document's rank.

Ranked query

A *ranked query* or *best match query*, as opposed to *exact match query* or *boolean query*, is a means of improving query results by assigning each result a weight according to its relevance for this particular query and displaying the results in descending order of relevance. Thus, the most relevant documents will be presented at the top of the result list, without having less relevant documents completely omitted as is the case with boolean queries, where documents not containing *all* of the query terms simply fall short. [1]
In so doing, *recall* gets immensely increased while some kind of *precision* is still achieved through the relevance ranking.

Yet still problematic is the concept of *relevance* [2] of a certain document for a particular query. The usual approach is to assign weights to the index terms and generate a document's weight (which acts as a measure for its relevance) from the matching term's weights.

lir.pl offers using various term weights calculated by different algorithms² [3] to allow for comparison of their effects and appropriateness as relevance measures.

Algorithms

Salton
Kascade einfach
Kascade komplex
Robertson
IDF
Custom
None

Termgewichte und Erschließung - Automatisches Klassifizieren

Verfahren

- Erstellen einer Kollektion mit intellektuell erschlossenen (klassifizierten) Dokumenten (bestehend aus Dokumenten in Klassen) = „Lernkollektion“
- Analyse der Termbeziehungen in den Dokumenten einer Klasse, z.B. auf der Basis einer Dokument-Term-Matrix der gewichteten Terme:

	Term1	Term2	Term3	Term4	Term5	Term6	Term7	Term8
Dok 1	0	4	0	0	0	2	1	3
Dok 2	3	1	4	3	1	2	0	1
Dok 3	3	0	0	0	3	0	3	0
Dok 4	0	1	0	3	0	0	2	0
Dok 5	2	2	2	3	1	4	0	2

- Ermittlung der häufigsten gemeinsamen (und hoch gewichteten) Terme einer Klasse
- Ermittlung der Häufigkeit dieser Terme in anderen Klassen
- Verwendung der geeigneten Terme als Repräsentationen der Klasse (idealerweise Terme, die für eine Klasse typisch sind, für andere nicht)

Ergebnis

Klassen und Terme, die die Klassen repräsentieren

Erstellen und Optimieren des Klassifizierers

- Festlegung der Bedingungen (Regeln), die zur Zuweisung eines Dokuments zu einer Klasse führen sollen, beispielsweise:
 - wie viele Terme einer Klasse müssen mindestens im Dokument enthalten sein
 - welche Gewichte müssen diese haben
- Testen des Klassifizierers anhand einer intellektuell erschlossenen Testkollektion
- ggf. Modifizierung der Regeln

Ergebnis

Klassifizierer, der algorithmisch die intellektuelle Erschließungspraxis der Testkollektion „simuliert“

Einsatz des Klassifizierers

- Anwendung des Klassifizierers für nicht erschlossene Dokumente: Zuteilung von Notationen eines Klassifikationssystems

Evaluierung des Klassifizierers

- idealerweise Evaluierung des Ergebnis durch Retrievaltests

Golub, Koraljka, Johan Hagelbäck und Anders Ardö: „Automatic classification using DDC on the swedish union catalogue“, in: Mayr, Philipp u. a. (Hrsg.): Preface of the 18th European Networked Knowledge Organization Systems Workshop (NKOS 2018), Porto: CEUR 2018 (CEUR workshop proceedings, vol 2020), S. 4–16.

Reiner, Ulrike: „Automatische DDC-Klassifizierung: bibliografische Titeldatensätze der Deutschen Nationalbibliografie“, in: Dialog mit Bibliotheken 22/1 (2010), S. 23–29.

Clustering

Ausgangspunkt

unstrukturierte, in der Regel sehr große Dokumentkollektion

Ziel

Strukturierung der Kollektion durch Ermittlung von Gruppen ähnlicher Dokumente

Verfahren

Berechnung der Ähnlichkeit von Dokumenten durch Analyse der Beziehungen zwischen Dokumenten und den in ihnen enthaltenen Termen und Festlegung eines Clustering-Algorithmus' für die Zuweisung von Dokumenten zu Clustern

Dokument-Term-Matrix,

d.h. welche Dokumente enthalten welche Terme mit welchem Gewicht

	Term1	Term2	Term3	Term4	Term5	Term6	Term7	Term8
Dok1	0	4	0	0	0	2	1	3
Dok2	3	1	4	3	1	2	0	1
Dok3	3	0	0	0	3	0	3	0
Dok4	0	1	0	3	0	0	2	0
Dok5	2	2	2	3	1	4	0	2

Dokument-Dokument-Matrix durch Berechnung der Skalarprodukte von jeweils zwei Dokumentvektoren

	Dok1	Dok2	Dok3	Dok4	Dok5
Dok1		11	3	6	22
Dok2	11		12	10	36
Dok3	3	12		6	9
Dok4	6	10	6		11
Dok5	22	36	9	11	

Erzeugen einer **Dokument-Beziehungs-Matrix** durch Festlegung eines Schwellenwertes (hier: 10)

	Dok1	Dok2	Dok3	Dok4	Dok5
Dok1		1	0	0	1
Dok2	1		1	1	1
Dok3	0	1		0	0
Dok4	0	1	0		1
Dok5	1	1	0	1	

Anwenden eines ‚**Clusteralgorithmus**‘ zur Verteilung der Dokumente auf Cluster

Cliquen-Algorithmus

alle Dokumente eines Clusters sind allen anderen Dokumenten des Clusters ähnlich; Dokumente in einem Cluster haben die engstmögliche Beziehung zueinander – Dokumente eines Clusters repräsentieren ein Thema (Topic)

Single-Link-Algorithmus

jedes Dokument eines Clusters ist mindestens einem Dokument des Clusters ähnlich; Dokumente eines Clusters haben schwache Beziehung zueinander – Dokumente eines Clusters repräsentieren keine Themen

Spielarten

(1) Verwendung von **Startclustern** und Berechnung von **Zentroiden**

- Festlegung von Clustern und beliebige Zuweisung von Dokumenten zu Clustern
- Berechnung eines Zentroids (d.h. eines Mittelwerts aller Dokumente eines Clusters)
- Berechnung der Ähnlichkeit zwischen den Dokumenten in den Clustern und den Zentroiden der Cluster und Neuverteilung der Dokumente in die Cluster
- Durchführung des Verfahrens bis zu stabilen Clustern

(2) **Hierarchisches Clustering**, z.B. durch

- iteratives Clustern von erzeugten Clustern bis hin zum einzelnen Dokument (Top-down)
- Berechnung von Zentroiden für die Cluster und Clustering der Zentroide (erzeugt erste hierarchisch höhere Ebene; Bottom-up)
- Fortführung des Prozesses bis zur gewünschten Hierarchie

Nutzen von Clustering im Information Retrieval

Termclustering

Clustering von Termen einer Kollektion erzeugt Mengen ähnlicher Begriffe, die für die automatische Erstellung thesaurus-ähnlicher Werkzeuge verwendet werden können:

- Ausweitung der Suche durch Einbeziehung ähnlicher Begriffe;
- Verlassen der strengen Matching-Bedingungen im Zeichenketten-Retrieval;
- Angleichung von Such- und Autorensprache;
- Visualisierung von Begriffsbeziehungen.

Dokumentclustering

Clustering von Dokumenten einer Kollektion erzeugt Mengen ähnlicher Dokumente, die für die Suche verwendet werden können:

- Ausweitung der Suche auf ähnliche Dokumente;
- Strukturierung von Treffermengen (NorthernLight-Prinzip);
- Visualisierung von Dokumentbeziehungen in Suchergebnissen;
- Verlassen der strengen Matching-Bedingungen im Zeichenketten-Retrieval;
- Relevance Feedback

Latent Semantic Indexing

Idee

- im **Vektorraummodell** sind Beziehungen zwischen Dokumenten und gewichteten Termen berechenbar:
 - ähnliche Dokumente** werden durch gemeinsame hochgewichtete Terme charakterisiert
 - ähnliche Terme** werden durch ihr gemeinsames Auftreten in verschiedenen Dokumente charakterisiert
- Cluster ähnlicher Terme entsprechen „**Konzepten**“
- Konzepte enthalten **Synonyme**, wodurch Indexierungen und Suchen mit Synonymen unterstützt werden

Deerwester, S, S Dumais und T Landauer: „Indexing by latent semantic analysis“, in: Journal of the American Society for Information Science 41/6 (1990), S. 391–407.

Titles:

- c1: *Human machine interface for Lab ABC computer applications*
- c2: *A survey of user opinion of computer system response time*
- c3: *The EPS user interface management system*
- c4: *System and human system engineering testing of EPS*
- c5: *Relation of user-perceived response time to error measurement*

- m1: *The generation of random, binary, unordered trees*
- m2: *The intersection graph of paths in trees*
- m3: *Graph minors IV: Widths of trees and well-quasi-ordering*
- m4: *Graph minors: A survey*

Terms	Documents								
	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

https://www.kde.cs.uni-kassel.de/wp-content/uploads/lehre/ws2006-07/IR/foalien/4Folie_02_IRmodels_d.pdf

Das Vektorraum-Modell

- Gehe davon aus, das t eindeutige Terme nach der Vorverarbeitung bleiben; nenne sie Indexterme oder das Vokabular.
- Diese “orthogonalen” Terme spannen einen Vektorraum mit Dimension t auf.
- Jedem Term i in einem Dokument oder einer Anfrage j wird ein reellwertiges Gewicht w_{ij} zugeordnet (im einfachsten Fall die Anzahl des Auftretens von i in j).
- Sowohl Dokumente als auch Anfragen werden als t -dimensionale Vektoren ausgedrückt:

$$d_j = (w_{1j}, w_{2j}, \dots, w_{tj})$$

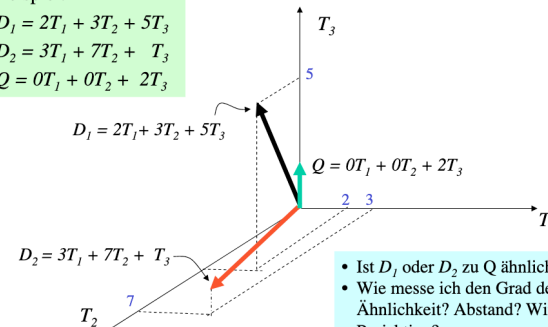
Grafische Darstellung

Beispiel:

$$D_1 = 2T_1 + 3T_2 + 5T_3$$

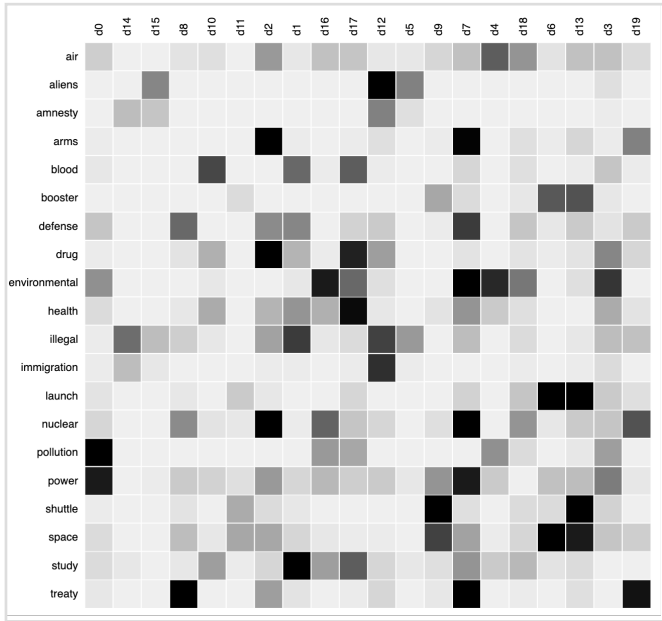
$$D_2 = 3T_1 + 7T_2 + T_3$$

$$Q = 0T_1 + 0T_2 + 2T_3$$



- Ist D_1 oder D_2 zu Q ähnlicher?
- Wie messe ich den Grad der Ähnlichkeit? Abstand? Winkel? Projektion?

Term-Dokument-Matrix

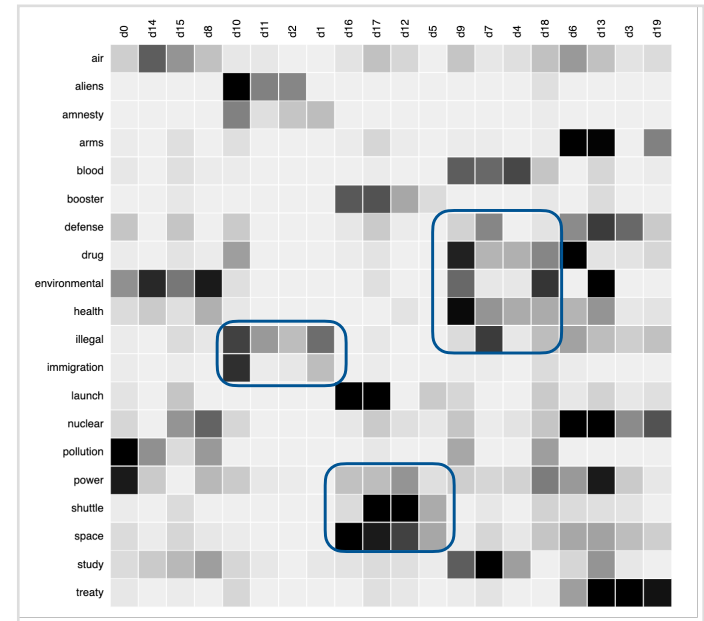


←
← Terme
←

↑ ↑ ↑
Dokumente

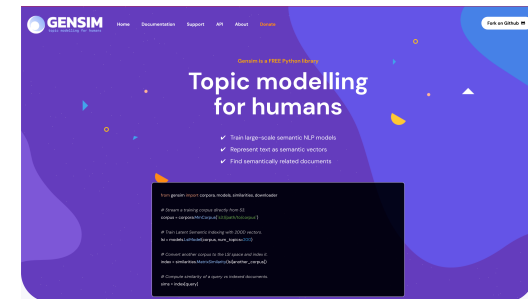
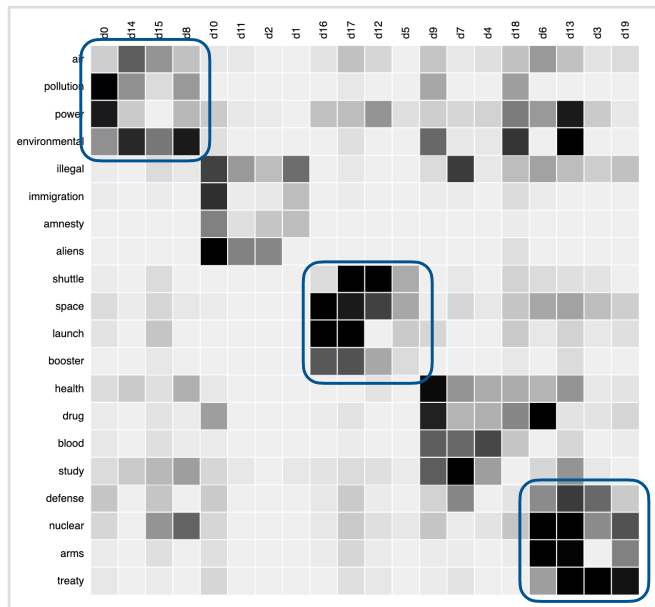
„ähnliche“ Terme
Konzepte

Ordnung nach Dokumentähnlichkeit



„ähnliche“ Dokumente

Ordnung nach Termähnlichkeit



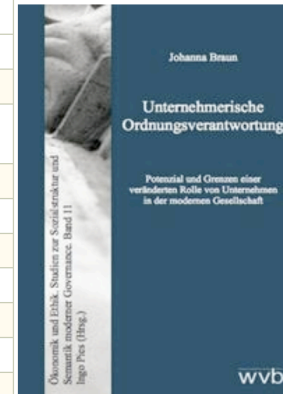
<https://radimrehurek.com/gensim/index.html>

Automatisches Erschließen

Automatische Zuteilung von normiertem Vokabular



Link zu diesem Datensatz	http://d-nb.info/1000028488
Art des Inhalts	Hochschulschrift
Titel	Unternehmerische Ordnungsverantwortung : Potenzial und Grenzen einer veränderten Rolle von Unternehmen in der modernen Gesellschaft / Johanna Braun
Person(en)	Braun, Johanna (Verfasser)
Verlag	Berlin : wvb. Wiss. Verl.
Zeitliche Einordnung	Erscheinungsdatum: 2009
Umfang/Format	XIX, 225 S. : graph. Darst. ; 23 cm, 410 gr.
Hochschulschrift	Zugl.: Halle (Saale), Univ., Diss., 2009
ISBN/Einband/Preis	978-3-86573-504-1 kart. : EUR 35.00
EAN	9783865735041
Sprache(n)	Deutsch (ger)
Beziehungen	Ökonomik und Ethik ; Bd. 11
Schlagwörter	Wirtschaftsordnung ; Wirtschaftsethik ; Unternehmensethik ; Corporate Social Responsibility
DDC-Notation	174.4 [DDC22ger]; 658.408 [DDC22ger]
Sachgruppe(n)	100 Philosophie ; 650 Management
Weiterführende Informationen	Inhaltsverzeichnis Kurzbeschreibung



Frankfurt	Signatur: 2010 A 24921 Bereitstellung in Frankfurt
Leipzig	Signatur: 2010 A 46327 Bereitstellung in Leipzig

Braun, Johanna: *Unternehmerische Ordnungsverantwortung : Potenzial und Grenzen einer veränderten Rolle von Unternehmen in der modernen Gesellschaft.*

Berlin: wvb. Wiss. Verl., 2009, XIX, 225 S.; graph. Darst
Zugl.: Halle (Saale), Univ., Diss., 2009
(Ökonomik und Ethik ; 11)

GND-Schlagwörter:

s Wirtschaftsordnung -- 9(DE-588)4066475-2; s Wirtschaftsethik -- 9(DE-588)4066439-9; s Unternehmensethik -- 9(DE-588)4202404-3; s Corporate Social Responsibility -- 9(DE-588)7697760-2

Maschinelle SLW: s Verantwortung; 040625478; 0.03378; s Unternehmen; 04061963X; 0.01504

VLB-Schlagwörter:

Ökonomische Ethik; Unternehmerische Verantwortung; Unternehmerische Praxis; Immaterielle Vermögenswerte; BC: Paperback; 1784: Hardcover, Softcover / Wirtschaft/Management

DDC: 174.4 -- c 174.4 -- eDDC22ger; 658.408 -- c 658.408 -- eDDC22ger

DNB-Sachgruppen: 100; 650; 100; 650

Munkelt, Johanna, Philipp Schaer und Klaus Lepsky: „Towards an IR test collection for the German National Library“, in: Gemulla, Rainer u. a. (Hrsg.): LWDA 2018: Lernen, Wissen, Daten, Analysen 2018; Proceedings of the Conference „Lernen, Wissen, Daten, Analysen“ Mannheim, Germany, August 22-24, 2018, Bd. 2191, Mannheim: CEUR-WS 2018 (CEUR-WS), S. 275–280.

<http://ceur-ws.org/Vol-2191/paper31.pdf>

Ressourcen für automatische Verfahren

Titelstichwörter

Verlagsschlagwörter

Scans von Inhaltsverzeichnissen

Braun, Johanna: *Unternehmerische Ordnungsverantwortung: Potenzial und Grenzen einer veränderten Rolle von Unternehmen in der modernen Gesellschaft.*
 Berlin: wvb. Wiss. Verl., 2009, XIX, 225 S.; graph. Darst
 Zugl.: Halle (Saale), Univ., Diss., 2009
 (Ökonomik und Ethik ; 11)

intellektuell vergebene Schlagwörter?

GND-Schlagwörter:

s Wirtschaftsordnung -- 9(DE-588)4066475-2; s Wirtschaftsethik -- 9(DE-588)4066439-9; s Unternehmensethik -- 9(DE-588)4202404-3; s Corporate Social Responsibility -- 9(DE-588)7697760-2

Maschinelle SLW: s Verantwortung; 040625478; 0.03378; s Unternehmen; 04061963X; 0.01504

VLB-Schlagwörter:

Ökonomische Ethik; Unternehmerische Verantwortung; Unternehmerische Praxis; Immaterielle Vermögenswerte;
 BC: Paperback; 1784: Hardcover, Softcover / Wirtschaft/Management

DDC: 174.4 -- c 174.4 -- eDDC22ger, 658.408 -- c 658.408 -- eDDC22ger

DNB-Sachgruppen: 100; 650; 100; 650

normiertes Vokabular mit
 Typ: „s“ = Sachschlagwort
 Schlagwort: Unternehmen
 Identnummer: 04061963X
 Konfidenzwert: 0.01504

Inhaltsverzeichnis	
Geleitwort des Herausgebers	VII
Vorwort	X
Inhaltsverzeichnis	XIII
Abbildungsverzeichnis	XV
Abkürzungsverzeichnis	XVII
Einführung	I
Kapitel 1: Unternehmensverantwortung im Spannungsfeld von Gewinn und Moral – eine kritische Würdigung paradigmatischer Positionen	12
1.1 Theoretische Ausgangspunkte	13
1.2 Primat der Verantwortung: Die Position Peter Ulrichs	20
1.2.1 Leitideen der „Integrativen Wirtschafts- und Unternehmensethik“	20
1.2.2 Kritische Würdigung	26
1.3 Zwischen Gewinnstreben und Verantwortung: Die Position Andreas G. Scherers	30
1.3.1 Leitideen der republikanischen Ethik der Multinationalen Unternehmung	30
1.3.2 Kritische Würdigung	36
1.4 Die Harmonie von Gewinnstreben und Verantwortung: Die Position Milton Friedmans	40
1.4.1 Leitideen der Friedmanschen Argumentation	40
1.4.2 Kritische Würdigung	44
1.5 Resümee des Kapitels	55
Kapitel 2: Unternehmerische Legitimität im Spannungsfeld von Marktwirtschaft und Demokratie – eine kritische Würdigung gesellschaftstheoretischer Konzeptionen	57
2.1 Unternehmerische Legitimität aus Sicht des deliberativen Republikanismus	58
2.1.1 Politik und Markt in der deliberativen Demokratie	59

XIV Inhaltsverzeichnis	
2.1.2 Die Habermas Rezeption von Peter Ulrich sowie von Andreas G. Scherer und Guido Palazzo	69
2.1.3 Kritische Würdigung	79
2.2 Unternehmerische Legitimität aus Sicht des ökonomischen Liberalismus	85
2.2.1 Politik und Markt im ökonomischen Liberalismus: Die Position Milton Friedmans	86
2.2.2 Kritische Würdigung	95
2.3 Unternehmerische Legitimität in marktwirtschaftlich verfassten Demokratien – eine ordnungsethische Perspektive	99
2.4 Resümee des Kapitels	113
Kapitel 3: Unternehmerische Ordnungsverantwortung – eine integrative Konzeption	115
3.1 Unternehmen als Träger von Verantwortung	116
3.2 Unternehmerische Ordnungsverantwortung im Kontext sozialer Dilemmata	125
3.3 Ordnungsverantwortung als Investition in unternehmerische Vermögenswerte	134
3.3.1 Investitionen in Humankapital	138
3.3.2 Investitionen in Organisationskapital	142
3.3.3 Investitionen in Reputationskapital	160
3.4 Instrumente strategischer Investitionsentscheidungen: Eine Illustration	170
3.5 Resümee des Kapitels	186
Zusammenfassung und Ausblick	189
Literaturverzeichnis	197
Personenregister	217
Sachregister	221

Mögliche Verfahrensweisen *

Linguistische Vorverarbeitung der Metadaten und Ressourcen

- OCR gescannter Ressourcen
- grammatikalische Normierung (Lemmatisierung), Dekomposition
- ggf. Entfernen unerwünschter Terme (Hochfrequenzterme)
- ggf. Gewichtung der Terme in Metadaten und Ressourcen

Linguistische Vorverarbeitung der normierten Terminologie

- ggf. Reduktion des Zeichensatzes und Entfernen störender Merkmale
- Auswertung der Synonymbeziehungen zur Synonymerkennung in Metadaten
- Auswertung der Sachgruppen für die Disambiguierung mehrdeutiger Terme
- Auswertung der hierarchischen Beziehungen für Aussagen zur Termspezifität

Abgleich von normierter Terminologie und Metadaten

- Erzeugung von potenziellen Indextermen
- Berechnung eines Qualitätsmaßes für zugeteilte Terme durch beispielsweise:
 - Berücksichtigung der Termgewichte der Terme in den Metadaten
 - Berücksichtigung der Sachgruppen der Terme
 - Berücksichtigung der Spezifität
 - Berücksichtigung bekannter Zuteilungswahrscheinlichkeiten aus erschlossenen Lernkollektionen
- Zuteilung aller Terme oberhalb eines festgelegten Wertes für die Zuteilungswahrscheinlichkeit

* Anm.: Dies ist keine Beschreibung des eingesetzten Algorithmus, da dessen genaue Arbeitsweise nicht veröffentlicht ist. Die Überlegungen basieren auf den vorhandenen bibliografischen Daten, den verfügbaren Erschließungsressourcen (Terminologie, Klassifikation) und den darin liegenden grundsätzlichen Möglichkeiten.

Kollektion

Fachdatenbank PHYS (inzw. Bestandteil von INSPEC) mit englisch-sprachiger Erschließung durch normiertes Vokabular (Deskriptoren) und Abstracts

Ziel von AIR/PHYS

automatische Indexierung der Dokumente mit Deskriptoren des PHYS-Thesaurus

Lernphase/Indexierungsidee

statistische Auswertung von 400.000 intellektuell erschlossenen Dokumenten, v.a. Untersuchung der Beziehung

Term \Rightarrow **z** \Rightarrow **Deskriptor**,

wobei **z** ein Maß für die Wahrscheinlichkeit ist, mit der ein **Deskriptor** einem Dokument (intellektuell) zugeteilt ist, wenn ein **Term** im Dokument vorhanden ist:

$$z = \frac{h(t,s)}{f(t)}$$

- $h(t,s)$ = Anzahl der Dokumente, in denen Term t vorkommt und Deskriptor s vergeben wurde
- $f(t)$ = Anzahl der Dokumente, in denen Term t vorkommt

Automatische Indexierung

- Aufbau eines „**Indexierungswörterbuchs**“ mit folgenden Elementen:
 - 350.000 **Term-Deskriptor-Beziehungen** mit $z > 0,3$
 - 70.000 **Synonymrelationen**
 - 200.000 **Deskriptor-Deskriptor-Beziehungen** mit $z > 0$ (als gewichtetes Maß für das gemeinsame Auftreten von Deskriptoren bei einem Dokument)
- **Automatische Indexierung** in zwei Phasen:
 - **Rohindexierung** mit regel- und lexikonbasierter Textanalyse und statistischer Relationierung
 - **Abgestimmte Indexierung** unter Einbeziehung von Deskriptor-Deskriptor-Relationen

Pilotanwendung AIR/PHYS

- Erschließung von 10.000 Dokumenten/Monat
- Zuteilung von im Schnitt **12 Deskriptoren je Dokument**
- **semi-automatisches Verfahren**: intellektuelle Nachbearbeitung mit ca. einem Drittel Korrekturbedarf

Retrievaltest (15.000 Dokumente, 300 Suchfragen)

- automatische Indexierung
Precision: 0.46
Recall: 0.57
- intellektuelle Indexierung
Precision: 0.53
Recall: 0.51

intellektuelle Bewertung der Erschließungsqualität durch Experten

- 1/3 intellektuelle Erschließung besser
- 1/3 automatische Indexierung besser
- 1/3 qualitativ gleichwertig

Knorz, Gerhard: „Automatische Indexierung“, Wissensrepräsentation und Information Retrieval, Potsdam: Universität 1994, S. 138–196.

Ziele

- **Anreicherung** von bibliografischen Referenzdaten aus dem Fach Jura
- Entwicklung einer **selektiven automatischen Indexierung** zur gewichteten Extraktion von Deskriptoren (SELIX)
- Entwicklung einer zuverlässigen Erkennung für **Themen-Aspekt-Beziehungen in Mehrwortgruppen** (THEAS)
- Durchführung eines umfangreichen **Retrievaltests**

Anreicherung der Titeldaten

- Scanning von Inhaltsverzeichnissen von ca. 3.000 Titeln aus dem Bestand Jura zur Verbreiterung der Indexierungsbasis (Abstracts: zu selten; Sachregister: problematisch)
- OCR mit newsWorks und MILOS-Rechtschreibkontrolle

Automatische Indexierung mit SELIX

- **Linguistische basierte Indexierung** mit
 - Grundformermittlung und Dekomposition für Sachtitel, Schlagwörter und Volltexte der Inhaltsverzeichnisse
 - Einbeziehung von semantischen Relationen (GND)
- **SELIX-Gewichtungsindexierung**

KASCADE - Gewichtung

Salton

$$\text{HfkImD}(g,d) * \log (n\text{Dok}(g)/n\text{AnzDok})$$

[Termhäufigkeit * log (Dokumenthäufigkeit / Dokumentenzahl), vgl. IDF]

Robertson

$$((K + 1) * \text{HfkImD}(g,d) / (K + \text{HfkImD}(g,d))) * \log((n\text{AnzDok} - n\text{Dok}(g) + 0.5) / (n\text{Dok}(g) + 0.5))$$

KASCADE

Kollektionsgewicht nG1

(nG1 ermittelt, ob eine Grundform für eine Dokumentenkollektion als Indexterm geeignet ist (für alle Dokumente (einer Kollektion) gleich))

$$nG1(g) = 1 - n\text{Dok}(g) / E(n\text{Dok}(g)) \quad (\text{für: } n\text{Dok}(g) < E(n\text{Dok}(g)); 0 \text{ sonst})$$

mit $E(n\text{Dok}(g)) = n\text{AnzDok} * (1 - \exp(-\lambda))$ wenn eine Poisson Zufallsverteilung angenommen wird:

$$P(i) = \exp(-\lambda) * (\lambda^i / i!) \quad \lambda = n\text{Coll}(g) / n\text{AnzDok}$$

Dokumentgewicht nG2

(nG2 ermittelt, ob eine Grundform für ein Dokument als Indexterm wichtig ist)

$$nG2(g,d) = (p(1) * 1 + \dots + p(\text{HfkImD}(g,d)) * \text{HfkImD}(g,d)) / \lambda$$

mit: $P(i) = \exp(-\lambda) * (\lambda^i / i!) \quad \lambda = n\text{Coll} * (n\text{DokLen} / n\text{CollLen})$

Längengewicht nG3

(nG3 bevorteilt längere Wörter im Gewichtungsverfahren (unabhängig von Dokument und Kollektion))

$$nG3(g) = \log (n\text{GruLen}(g)) / 4$$

Gewichtungsfunktion

$$nG = F1 * nG1 + F2 * nG2 + F3 * nG3$$

(wobei F1-F3 frei wählbar sind (Standard: 1))

Hüther, Hubert: „Selix im DFG-Projekt Kascade“, in: Zimmermann, Harald H. und Volker Schramm (Hrsg.): Knowledge management und kommunikationssysteme, workflow management, multimedia, knowledge transfer - proceedings des 6. Internationalen symposiums für informationswissenschaft (ISI '98), prag 3.-7. November 1998, Bd. 34, Hochschulverband für Informationswissenschaft 1998 (Schriften zur informationswissenschaft), S. 397–403.

KASCADE - Retrievaltest

Rahmenbedingungen

- 3.000 Referenzdatensätze aus dem Fach Jura
- alle angereichert um Inhaltsverzeichnisse im Volltext
- 60 von Juristen formulierte Suchthemen
- Testdurchführung durch Projektmitarbeiter
- Relevanzbewertung durch Juristen
- Recall-Berechnung nach Pooling-Methode

Besonderheiten bei den Suchthemen

- breite thematische Streuung – speziell neben allgemein
- viele Komposita und Mehrwortbegriffe
- viele komplexe Themen, d.h. Themenverknüpfungen
- nur 15% Einwort-Suchthemen (mit nur einem Nichtkompositum)

Ergebnisse

	Mittelwerte von		Null-Treffer-Suchen
	Recall	Precision	
Titel und Deskriptor (automatisch indexiert)	0.06	0.98	42
Titel, Deskriptor, Inhaltsverz. (nicht automatisch indexiert)	0.54	0.75	7
Titel, Deskriptor, Inhaltsverz. (automatisch indexiert)	0.92	0.70	4

Weiterführende Literatur

Literatur zur Informationserschließung

Automatisches Indexieren

Automatisches Klassifizieren

Texttechnologie

Allen, James: Natural language understanding, Redwood City, Calif. [u.a.]: Benjamin/Cummings Publ. 1995.

Carstensen, Kai-Uwe: Computerlinguistik und Sprachtechnologie : eine Einführung, Heidelberg: Spektrum, Akad. Verl. 2010.

Gödert, Winfried, Klaus Lepsky und Matthias Nagelschmidt: Informationserschließung und Automatisches Indexieren : ein Lehr- und Arbeitsbuch, Berlin [u.a.]: Springer 2012 (X.media.press).

Hausser, Roland: Grundlagen der Computerlinguistik : Mensch-Maschine-Kommunikation in natürlicher Sprache, Berlin [u.a.]: Springer 2000.

Jurafsky, Daniel und James H. Martin: Speech and Language Processing : an introduction to natural language processing, computational linguistics, and speech recognition, Pearson Prentice Hall/Pearson education international 2009 (Pearson international edition).

Lobin, Henning: Computerlinguistik und Texttechnologie, Paderborn: W. Fink 2010.

Lepsky, Klaus: „Automatische Indexierung“, in: Kuhlen, Rainer, Wolfgang Semar und Dietmar Strauch (Hrsg.): Grundlagen der praktischen Information und Dokumentation: Handbuch zur Einführung in die Informationswissenschaft und -praxis, Berlin: De Gruyter 2013, S. 272–285.

Manning, Christopher und Hinrich Schütze: Foundations of statistical natural language processing, MIT Press 1999.