

Ulrike Reiner

Automatische DDC-Klassifizierung

Bibliografische Titeldatensätze der Deutschen Nationalbibliografie

Das Klassifizieren von Objekten (z. B. Fauna, Flora, Texte) ist ein Verfahren, das auf menschlicher Intelligenz basiert. In der Informatik – insbesondere auf dem Gebiet der Künstlichen Intelligenz (KI) – wird u. a. untersucht, inwieweit Verfahren, die menschliche Intelligenz benötigen, automatisiert werden können. Hierbei hat sich herausgestellt, dass die Lösung von Alltagsproblemen eine größere Herausforderung darstellt, als die Lösung von Spezialproblemen, wie z. B. das Erstellen eines Schachcomputers. So ist »Rybka« der seit Juni 2007 amtierende Computerschach-Weltmeister. Inwieweit Alltagsprobleme mit Methoden der KI gelöst werden können, ist eine – für den allgemeinen Fall – noch offene Frage. Beim Lösen von Alltagsproblemen spielt die Verarbeitung der natürlichen Sprache, wie z. B. das Verstehen, eine wesentliche Rolle. Den »gesunden Menschenverstand« als Maschine (in der Cyc-Wissensbasis in Form von Fakten und Regeln) zu realisieren, ist Lenat's Ziel seit 1984. Bezüglich des KI-Paradeprojektes »Cyc« gibt es Cyc-Optimisten und Cyc-Pessimisten.

Das Verstehen der natürlichen Sprache (z. B. Werktitel, Zusammenfassung, Vorwort, Inhalt) ist auch beim intellektuellen Klassifizieren von bibliografischen Titeldatensätzen oder Netzpublikationen notwendig, um diese Textobjekte korrekt klassifizieren zu können. Seit dem Jahr 2007 werden von der Deutschen Nationalbibliothek (DNB) nahezu alle Veröffentlichungen mit der Dewey-Dezimalklassifikation (DDC) intellektuell klassifiziert. Die Menge der zu klassifizierenden Veröffentlichungen steigt spätestens seit der Existenz des World Wide Web schneller an, als sie intellektuell sachlich erschlossen werden kann. Daher werden Verfahren gesucht, um die Klassifizierung von Textobjekten zu automatisieren oder die intellektuelle Klassifizierung zumindest zu unterstützen. Seit 1968¹⁾ gibt es Verfahren zur automatischen Dokumentenklassifizierung (Information Retrieval, kurz: IR) und seit

1992²⁾ zur automatischen Textklassifizierung (Automated Text Categorization, kurz: ATC). Seit immer mehr digitale Objekte im World Wide Web zur Verfügung stehen, haben Arbeiten zur automatischen Textklassifizierung seit ca. 1998 verstärkt zugenommen. Dazu gehören seit dem Jahr 1996 auch Arbeiten zur automatischen DDC-Klassifizierung bzw. RVK (Regensburger Verbundklassifikation)-Klassifizierung von bibliografischen Titeldatensätzen und Volltextdokumenten. Bei den Entwicklungen handelt es sich unseres Wissens bislang um experimentelle und keine im ständigen Betrieb befindlichen Systeme. Auch das Projekt der Verbundzentrale des Gemeinsamen Bibliotheksverbundes (VZG-Projekt) Colibri/DDC ist seit 2006 u. a. mit der automatischen DDC-Klassifizierung befasst. Die diesbezüglichen Untersuchungen und Entwicklungen dienen zur Beantwortung der Forschungsfrage: »Ist es möglich, eine inhaltlich stimmige DDC-Titelklassifikation aller GVK-PLUS (Gemeinsamer Verbundkatalog (GVK) und Online Contents (OLC))-Titeldatensätze automatisch zu erzielen?«

Colibri/DDC-Wettbewerb

Da in der Fachwelt starkes Interesse an (semi-)automatischen Klassifizierungssystemen von Textobjekten besteht und um den Anreiz für die Entwicklung solcher Systeme zu steigern, wurde auf dem 98. Deutschen Bibliothekartag in Erfurt der »Colibri/DDC-Wettbewerb«³⁾ initiiert, mit dem Ziel, den besten automatischen DDC-Klassifizierer für bibliografische Titeldatensätze zu finden. In Abb. 1 ist ein fiktives Szenario für den Colibri/DDC-Wettbewerb skizziert: Während sich die Klassifizierungskomponente `vc_dcl` (`vzg_colibri_ddc_classifier`) des DDC-Suchsystems `vc_ds` ((Reiner 2009b), S. 18 – 19 und S. 39) und das System Y schon länger in der Entwicklung befinden, gibt es ein System X, das kurz vor dem Start steht und ein System Z, das hiervon noch etwas

Lösung von Alltagsproblemen durch Künstliche Intelligenz?

Colibri-Projekt

Suche nach Verfahren zur automatischen Klassifizierung

entfernt ist. Unabhängig vom Startzeitpunkt kann sich der Weg zum Erfolg unterschiedlich gestalten: Die Entwicklung der Systeme X und Y ist unbekannt; das zuletzt an den Start gegangene System Z könnte alle anderen Systeme »überholen«, dann jedoch bei Verbesserungsversuchen in Stagnation geraten, ohne ans Ziel zu gelangen. vc_dcl könnte innerhalb der Entwicklung Rückschläge erfahren und überflüssige (Rück-)Wege beschreiten, dann jedoch (wie erhofft) ans Ziel gelangen. Ob das illustrierte, fiktive Szenario des Colibri/DDC-Wettbewerbs in dieser oder ähnlicher Weise eintreten und ein einsatzbereiter, automatischer DDC-Klassifizierer Realität werden wird, wird die Zukunft zeigen.

Ziel:
Automatischer
DDC-Klassifizierer

Nach (Salton/McGill 1983)⁵⁾ sind für einen Systemtest mindestens folgende Bestandteile notwendig:

- Modell des Systems oder detaillierte Beschreibung des Systems und seiner Komponenten,
- zu testende Hypothesen,
- Bewertungskriterien und Maße, die diese Kriterien widerspiegeln und
- Methoden, die Daten zu ermitteln und zu bewerten.

Notwendige
Bestandteile

Nach diesen vier Kriterien wird nachfolgend die Colibri/DDC-Klassifizierungskomponente vc_dcl genauer beschrieben und es werden zusätzlich einige andere in der Entwicklung befindlichen Systeme zur automatischen Klassifizierung zum Vergleich herangezogen.

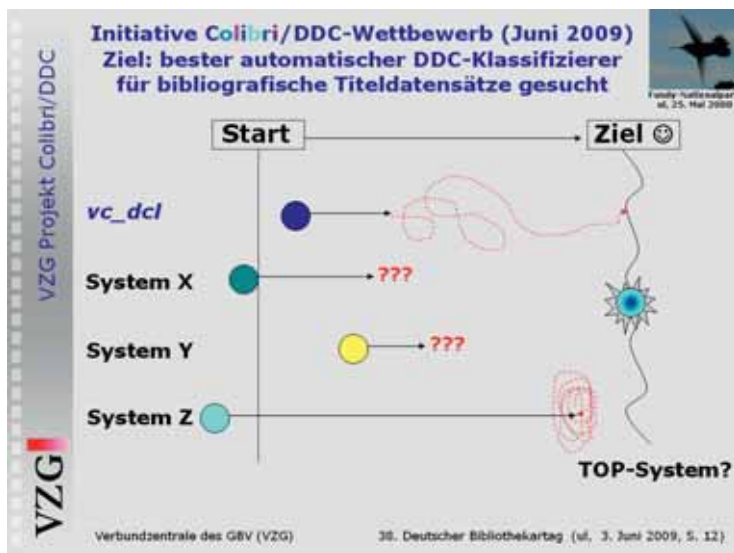


Abb. 1: Initiative Colibri/DDC-Wettbewerb (Juni 2009)

Systeme zur automatischen Klassifizierung

Mit der Entwicklung der automatischen Klassifizierungskomponente vc_dcl des Suchsystems vc_ds innerhalb des Projektes Colibri/DDC⁶⁾ wurde im Jahr 2006 begonnen und sie wurde im Jahr 2008 fortgeführt. Die klassifizierte Textobjekte sind bibliografische Titeldatensätze der Deutschen Nationalbibliografie. Da die Titeldatensätze intellektuell vergebene DDC-Notationen enthalten, können daraus DDC-Testbestände im IR und in der KI auch Menge der Testdokumente bzw. Testkorpus genannt) erstellt und (durch den Vergleich der intellektuell vergebenen mit den automatisch ermittelten DDC-Notationen) DDC-Klassifizierungssysteme bewertet werden (s. u. Bewertungskriterien und -maße). Ein Aspekt des Colibri/DDC-Wettbewerbs ist - analog der Tradition des Information Retrieval (Cranfield, TREC, GOV2, CLEF, REUTERS u. a.) - Standard-Testbestände für die automatische DDC-Klassifizierung aufzubauen/festzulegen. Bislang stehen die DNB-DDC-Testbestände in_DNB-2007, in_DNB-2009 und in_DNB-2009-2 im Colibri/DDC-Wettbewerb zur Verfügung, s. Abb. 2. Für die automatische Klassifizierung werden die MAB2-Dateien als Eingabedateien in_DNB-2007, in_DNB-2009 und in_DNB-2009-2 aufbereitet und als Menge von Objekt-Attribut-Wert-Tripeln (DDC-Notation, Deskriptor, Deskriptorwert)⁷⁾ repräsentiert. Inkor-

Bibliografische
Titeldatensätze
als Textobjekte

Systemtests

Für einen Qualitätsvergleich von Systemen sind vergleichbare Testbedingungen unerlässlich. Diese sind zurzeit bei automatischen Klassifizierungssystemen noch nicht gegeben. Die Systeme differieren z. B. in Annahmen, Voraussetzungen, Anwendungsbereichen, zu klassifizierenden Objekten, Testbeständen und Bewertungsmaßen (ein Zustand, wie er 1983 im Information Retrieval vorherrschte⁸⁾). Um für den »besten automatischen DDC-Klassifizierer« vergleichbare Testbedingungen zu erstellen, wurde zur Standardisierung der Colibri/DDC-Wettbewerb ins Leben gerufen.

rekte DDC-Notationen und irrelevante Deskriptorwerte werden vorher eliminiert und relevante Deskriptorwerte u. a. transliteriert, deren Sonderzeichen entfernt und in Kleinschreibung transformiert.

DNB-DDC-Testbestände	Anzahl Titeldatensätze ¹⁾ klassifizierte	Deskriptorwerte ²⁾	Sprache			
			Deutsch	Englisch	andere	
in_DNB-2007	25.653	15.365	10,5	76%	20%	4%
in_DNB-2009	30.717	21.422	11,0	80%	18%	2%
in_DNB-2009-2	45.935	33.536	10,5	78%	20%	2%

Abb. 2: Charakteristika dreier DNB-DDC-Testbestände

Die Klassifizierungsbasis für die DDC-Testbestände bildet GVK-DDC (die Teilmenge des Gemeinsamen Verbundkataloges GVK, deren GVK-Titeldatensätze mindestens eine DDC-Notation enthalten), die sich durch DDC-Lieferanten wie z. B. DNB, LoC und OCLC jährlich vergrößert:

GVK-DDC	Anz. Titeldatensätze	GVK-DDC 0	GVK-DDC 1	GVK-DDC 2	GVK-DDC 3	GVK-DDC 4	GVK-DDC 5	GVK-DDC 6	GVK-DDC 7	GVK-DDC 8	GVK-DDC 9
2004	3,0 Mio.	99M	57M	109M	564M	46M	145M	306M	188M	324M	280M
2008	4,3 Mio.	236M	153M	253M	1.3G	118M	386M	677M	390M	666M	571M
2009	5,9 Mio.	515M	271M	394M	2.2G	358M	714M	1.2G	629M	981M	858M

Abb. 3: GVK-DDC mit Teilmengen GVK-DDC0 ... GVK-DDC9 (2004 - 2009)

Aus GVK-DDC wird die DDC-Datenbasis vc_DB (DataBase) und aus dem DDC-System die DDC-Wissensbasis vc_KB (Knowledge Base) gebildet. vc_DB und vc_KB bilden zusammen

vc_DB_PLUS, aus Effizienzgründen invertiert (vc_IDB: Inverted DataBase). Die Einlesezeit von vc_IDB-2004 (510 MB) in den Hauptspeicher beträgt drei Min., von vc_IDB-2008 (712 MB) fünf Min. und von vc_IDB-2009 (925 MB) sechs Min. Die DNB-DDC-Testbestände in_DNB-2007, in_DNB-2009 und in_DNB-2009-2 werden größenabhängig in zwei bis drei Stunden automatisch klassifiziert. Der Klassifizierungsalgorithmus verwendet IR- und KI-Verfahren ((Reiner 2009a), S. 7 - 13; (Reiner 2009b), S. 15 - 17), jedoch bislang keine linguistischen Verfahren. Das Ergebnis der automatischen Klassifizierung wird intellektuell und automatisch bewertet. Weitere Klassifizierungsprojekte s. Abb. 4.

Keine Anwendung linguistischer Verfahren

Hypothesen

Einzelne Hypothesen, Annahmen etc. werden hier exemplarisch wiedergegeben.⁸⁾

Auto-DDC (Forschungsfragen):

- Sind bibliothekarische Klassifikationsschemata für eine maschinelle Verarbeitung geeignet?
- Sind Methoden des überwachten Lernens geeignet, um effektiv auf bibliografischen Daten zu arbeiten?

Klassifizierungsprojekt	Auto-DDC	Auto Dewey	Colibri/DDC	Pfeffer/RVK	Scorpion/DDC	TopicModels/DDC
Laufzeit	seit 2006	seit 2006	seit 2006/2008	seit 2005	1996-2000	seit 2009
Klassifizierung	DDC	DDC	DDC	RVK	DDC	DDC
zu klassifizierende Textobjekte	bibliograf. Titeldatensätze	bibliograf. Titeldatensätze	bibliograf. Titeldatensätze	bibliograf. Titeldatensätze	Elektron. Webdokumente (Volltexte)	bibliograf. Titeldatensätze
berücksichtigte Daten der zu klassifizierenden Textobjekten	inhaltsbezogene Kategorien (MARC21-Titeldatensätze)	engl., franz., ital. literar. Autoren: Dichtung, Dramen, Prosa	Deskriptorwerte von inhaltstragenden Deskriptoren (MAB2-Kategorien)	Titel-/Schlagwörter	„Editorial Support System (ESS)“-Datensätze des DDC-Systems	OAI Metadaten: Titel-/Schlagwörter, Zusammenfassungen
Testbestände	22.110 LoC-Titeldatensätze aus Science & Technology (BDS&T)	unbekannt	in-DNB-2007 in-DNB-2009 in-DNB-2009-2	10.000 zufällige Titeldatensätze aus Fallbasis	NetFirst ([Shafer 1997], s. Fußnote 34)	719 englische Titeldatensätze 1.000 deutsche Titeldatensätze
Basis für die Klassifizierung	66.440 LoC-Titeldatensätze aus Science & Technology (BDS&T)	MARC21-Titeldatensätze	GVK-DDC-Titeldatensätze & DDC-System (vc_IDB)	SWB-/HeBIS-Titeldatensätze (Fallbasis)	DDC-System als Wissensbank (Scorpion Dewey Database)	100.000 BASE1-Dokumente, Wikipedia (Trainingsdokumente)

Algorithmen aus dem/den Gebiet/en	KI-ATC: Naïve Bayes, Support Vector Machines, DDC restructuring	unbekannt (Menge von Algorithmen)	IR: Ähnlichkeitsmaß & KI: heuristische Funktion	KI: fallbasiertes Schliessen	IR: Ähnlichkeitsmaße (SMART ²)	KI-ATC: Support Vector Machines, Latent Semantic Analysis, SEQ-based classifier, Wikipedia-based classifier
Hypothesen / Fragen	ja	unbekannt	ja	nein	ja	ja
Bewertung	automatisch	unbekannt	intellektuell u. automatisch	intellekt. u. automatisch	automatisch	automatisch
Automatisierungsgrad	semi-automatisch	(semi-) automatisch	automatisch	automatisch	automatisch	automatisch

¹ [Pieper/Summann 2006] Dirk Pieper; Friedrich Summann: Bielefeld Academic Search Engine (BASE): An end-user oriented institutional repository search service. Library Hi Tech, Bd. 24, Nr. 4, 2006, pp. 614-619. http://eprints.rclis.org/9160/1/pieper_summann_final_web.pdf

² [http://en.wikipedia.org/wiki/SMART_Information_Retrieval_System]

Abb. 4: Klassifizierungsprojekte als potenzielle Kandidaten im Colibri/DDC-Wettbewerb⁹⁾

Colibri/DDC (Hypothesen)

- Es gibt signifikante Unterschiede zwischen
 1. unterschiedlichen DDC-Daten-/Wissensbasen,
 2. den DDC-Klassen, 3. den DNB-Reihen A, B und H und 4. hinsichtlich der Stelligkeit der DDC-Notationen.
- Es gibt keine signifikanten Unterschiede zwischen 5. unterschiedlichen Testbeständen und 6. englischen und deutschen DNB-Titeldatensätzen.

Pfeffer/RVK (Testbeschränkungen)

- keine Reihen, keine Zeitschriften, keine formalen Klassifikationen.

Scorpion/DDC (Forschungsfragen)

- Welchen Effekt hätte die zusätzliche Verwendung von »Library of Congress Subject Headings« (LCSH)?
- In welcher Weise änderten sich die Ergebnisse, wenn andere Algorithmen als ATN/ATC verwendet würden?
- Welches Ergebnis lieferten andere Klassifikationen, wie z. B. die »Library of Congress Classification« (LCC) unter ähnlichen Bedingungen?

TopicModels/DDC (Hypothese)

- Der Inhalt eines wissenschaftlichen Dokumentes wird mithilfe seines Titels, Schlagwörtern und einer kurzen Beschreibung zuverlässig klassifiziert.

Bewertungskriterien und -maße

Um die Qualität der Klassifizierungsergebnisse beurteilen zu können, müssen diese bewertet und wegen der Vergleichbarkeit müssen dieselben Bewertungskriterien und -maße verwendet werden. Dies ist zurzeit bei der automatischen Klassifizierung von bibliografischen Titeldatensätzen noch nicht der Fall. Bei den Klassifizierungsprojekten sind verschiedene Bewertungsmaße im Einsatz. Außerdem kann es sein, dass Berechnungen mit dem selbem Maß unterschiedlich vorgenommen werden (Micro/Macroaverage von Precision/Recall⁹⁾). Im Projekt Colibri/DDC ist durch die DNB eine intellektuelle (exakte, gute, mittlere, keine, ausreichende Übereinstimmung, Hauptstichgruppentreffer) und durch die vc_ds-Programmkomponente vc_dce (vzg colibri_ddc classification results evsaluator) eine automatische Bewertung (u. a. mit den Bewertungsmaßen C, CP und CN) und anhand von Beispielen eine vergleichende Betrachtung von Bewertungen der Klassifizierungsergebnisse vorgenommen worden ((Reiner 2009a), S. 13 ff.; (Reiner 2009b), S. 31 ff.).

Zurzeit sind verschiedene Bewertungsmaße im Einsatz

Andere Projekte verwenden:

- Auto-DDC: Macro-/(depth-based)Microaverage von Precision/Recall;
- Pfeffer/RVK: perfekte, gute, mäßige und schlechte Übereinstimmung;

- Scorpion/DDC: Relationen: übereinstimmend, allgemeiner als, korreliert, synonym, exakt und thematisch nah;
- TopicModels/DDC: Precision, Recall und F-Score.

Tests und Bewertung

In den o. g. Projekten sind einige Tests mit unterschiedlichen Bewertungen durchgeführt worden. Diese Arbeit soll die notwendige Reproduzierbarkeit der Tests und die Vergleichbarkeit unterschiedlicher Klassifizierungsergebnisse fördern, um den besten automatischen (DDC-)Klassifizierer für bibliografische Titeldatensätze ermitteln zu können.

Ergebnisse und Ausblick

Für den Stand der Entwicklungen seien einzelne Ergebnisse/Einschätzungen wiedergegeben:

- Scorpion/DDC: »Scorpion cannot replace human cataloging. There are many aspects of

human cataloging that are difficult if not impossible to automate.« (Shafer 1996)

- LCC: »Fully automatic classification might not be possible considering the size and diversity of the LCC scheme« (Larson 1992)¹⁰.
- Auto-DDC: »With no more than three interactions, a classification accuracy of nearly 90 % is achieved, thus providing a practical solution to the automatic bibliographic classification problem« (Wang 2009), S. 2269.
- Pfeffer/RVK: Die SWB-Datenbasis mit 2.496.839 Titeln erzielte 57,26 % (Hamming) und 56,89 % (IDF) »perfekte« und 18,99 % (Hamming) und 18,84 % (IDF) »gute« Ergebnisse (Pfeffer 2008), S. 13; die intellektuelle Überprüfung (SWB, HeBIS) kommt jedoch zu dem Schluss, »dass die Qualität der automatisch generierten Klassifikationen zu schlecht ist, um direkt in die Verbunddatenbanken eingespielt zu werden.«¹¹
- TopicModels/DDC: »By this procedure we get a classification value for each main class of the DDC which expresses the relatedness of a given OAI-input stream to the selected class ... With an



Stadtbibliothek Tübingen (Deutschland) · www.kraas-lachmann.com

Als wär man da.

Ihre Nutzer wollen bequem am Strand auf die Inhalte Ihrer wertvollen Originalausgaben zugreifen? Kein Problem! Wir beherrschen mit unseren Digital- und Analogsystemen alle Prozesse der Dokumenten-Erfassung, -Archivierung, -Verarbeitung und -Bereitstellung. Seit mehr als 40 Jahren.

Zeutschel, die Zukunft der Vergangenheit.



Zeutschel GmbH · Heerweg 2 · 72070 Tübingen · Tel.: +49 7071 9706-0
Fax: +49 7071 9706-44 · info@zeutschel.de · www.zeutschel.de

overall F-score of .761, SVMs provide an adequate DDC-related method to classify documents based on their OAI metadata.« (Mehler/Waltinger 2009), S. 8/14.

- Colibri/DDC: In 63.85 % (CN-Wert) stimmen intellektuell vergebene und automatisch ermittelte DDC-Notationen in der DDC-Hauptklasse überein (bei vc_DB-2008 und in_DNB-2009 unter genannten Voraussetzungen in (Reiner 2009a), Fußnoten 45 und 52). Weitere Ergebnisse in (Reiner 2009b), S. 41 - 47. Durch Eliminierung weiterer Banalwörter erhöht sich der CN-Wert um ca. 1 %. Der neueste Testbestand in_DNB-2009-2 liefert mit der neuesten Klassifizierungsbasis GVK-DDC-2009 einen CN-Wert von ca. 91 %! Dieses erstaunliche Ergebnis muss näher untersucht werden (von insgesamt 45.935 Titeldatensätzen werden nur 13.115 Titeldatensätze klassifiziert). Mit hoher Wahrscheinlichkeit ist eine große Menge weiterer Titeldatensätze des DNB-Testbestandes in_DNB-2009-2 in der neuen Klassifizierungsbasis GVK-DDC-2009 enthalten, ohne mit den derzeitigen DNB-Kennungs-, ISBN- und ISSN-Prüfungen erkannt worden zu sein. Wenn dies der Fall ist, müssen weitere Titeldatensatz-Identifizierungs-Prüfungen (auf Enthaltensein in der Klassifizierungsbasis), z. B. auf Titel-Person-Basis, implementiert werden.

Fazit

Die inspizierten automatischen (DDC-)Klassifizierer arbeiten besser als der Zufall, aber für einen professionellen umfangreichen Einsatz für Mio. von

zu klassifizierenden (z. B. GVK-)Titeldatensätzen sind sie noch nicht geeignet. Auto-DDC erreicht zwar 90 % »accuracy«, dies wird jedoch nur in Kombination mit der intellektuellen Leistung (bis zu drei Benutzeraktionen) erreicht. In den nächsten Schritten müssen die Klassifizierungsergebnisse verbessert und wegen der angestrebten Vergleichbarkeit im Colibri/DDC-Wettbewerb dieselben Hypothesen geprüft und die Ergebnisse in gleicher Weise bewertet werden.

Dank

Die Autorin dankt Kristina Knull-Schlomann (DNB) für Ihre Ermunterung, hier eine schriftliche Version¹²⁾ ihres Vortrages auf dem 98. Deutschen Bibliothekartag in Erfurt in aktualisierter Form zu veröffentlichen und den DNB- und VZG-KollegInnen für die zur Verfügungstellung der Daten zur Erstellung der DNB-DDC-Testbestände und GVK-DDC-Klassifizierungsbasen, insbesondere (in alphabetischer Reihenfolge): Angelika Cremer-Reiber (DNB), Reiner Diedrichs (VZG), Siegfried Kalb (VZG) und Claudia Werner (DNB).

Anschrift von Dr. Ulrike Reiner: Projektleitung Colibri / DDC
Verbundzentrale des Gemeinsamen Bibliotheksverbundes (VZG),
Platz der Göttinger Sieben 1, 37073 Göttingen,
E-Mail: ulrike.reiner@gbv.de

Automatische
Klassifizierer
arbeiten nicht
professionell
genug

Anmerkungen

- 1 Salton, Gerard: Automatic Information Organization and Retrieval. McGraw-Hill, Inc., New York u. a., 1968, S. 112, 133 - 135; Bollmann, Peter; Konrad, Erhard; Schneider, Hans-Jochen; Zuse, Horst: Anwendung automatischer Klassifikationsverfahren mit dem System FAKYR. In: Wolfgang Dahlberg (Hrsg.): Kooperation in der Klassifikation. Proceedings der Sekt. 1 - 3 der 2. Fachtagung der Gesellschaft für Klassifikation e. V., Frankfurt-Höchst, 06. - 07. April 1978. Frankfurt, Gesellschaft für Klassifikation, 1978, S. 156 - 165.
- 2 Lewis, David D.: An Evaluation of Phrasal and Clustered Representations on a Text Categorization Task. Proceedings of SIGIR-92, 15th ACM International Conference on Research and Development in Information Retrieval, ACM Press, New York, US, 1992, pp. 37 - 50.

- 3 (Reiner 2009b) Reiner, Ulrike: Automatische DDC-Klassifizierung von bibliografischen Titeldatensätzen. 98. Deutscher Bibliothekartag: Ein neuer Blick auf Bibliotheken. TK10: Information erschließen und recherchieren. Blockveranstaltung: Inhalte erschließen – mit neuen Tools, Erfurt, 3.6.09, S. 12.
<<http://www.opus-bayern.de/bib-info/volltexte/2009/736/>>
- 4 Reiner, Ulrike: Experimente im Gebiet des Information Retrieval – Überblick und Stand der Forschung. Technische Universität Berlin, Institut für Angewandte Informatik, Fachbereich Informatik. LIVE-Bericht Nr. 6/83 (entstanden im Rahmen des vom BMFT geförderten Projektes »Leistungsbewertung von Information Retrieval Verfahren« (LIVE)), 1983.
<<http://portal.acm.org/citation.cfm?id=253495.253527>>
- 5 Salton, Gerard; McGill, Michael J.: Introduction to Modern Information Retrieval. Chapter 5: Retrieval Evaluation. McGraw-Hill International Book Company, Hamburg u. a., 1983, S. 158.
- 6 (Reiner 2009a) Reiner, Ulrike: VZG-Projekt Colibri. Bewertung von automatisch DDC-klassifizierten Titeldatensätzen der Deutschen Nationalbibliothek (DNB). VZG-Colibri-Bericht 1/2008, August 2008 – Februar 2009.
<<http://taipan.dyndns.org/~ul/colibri05.pdf>>
- 7 DDC-Notation: Notation einer DDC-Klasse; Deskriptor: Pica+-Kategorie bzw. MAB2-Kategorie, deren Werte zur inhaltlichen Charakterisierung beitragen. Zurzeit berücksichtigte Pica+- bzw. MAB2-Deskriptoren s. (Reiner 2009b), S. 24 und 25; Deskriptorwert: Wert eines Deskriptors, s. (Reiner 2009b), S. 14.
- 8 Auto-DDC (Wang 2009) Wang, Jun: An Extensive Study on Automated Dewey Decimal Classification. Journal of the American Society for Information Science and Technology, Vol. 60, No. 11, 2009, pp. 2269 – 2286.
AutoDewey Tillett, Barbara B.: Library of Congress Report. ALA ALCTS Committee on Cataloging: Description and Access. Midwinter Meeting, Philadelphia, PA, January 12, 2008.
<<http://www.libraries.psu.edu/tas/jca/ccda/docs/lc0801.pdf>>
Green, Rebecca: Literary Authors: AutoDewey and LC Name Authority File. Dewey Breakfast/Update. ALA Midwinter Meeting, January 12, OCLC, 2008.
<http://www.oclc.org/dewey/discussion/papers/literary_authors.ppt>
Colibri/DDC (Reiner 2009a); (Reiner 2009b).
Pfeffer/RVK: <http://blog.bib.uni-mannheim.de/Classification/wp-content/uploads/2008/6/bibtag2008_rvk.pdf>
Scorpion/DDC: <<http://www.worldcat.org/oclc/54388055>>
<<http://www.worldcat.org/oclc/54084278>>
<<http://worldcat.org/arcviewer/1/OCC/2003/06/03/0000003411/viewer/file1.html>>
<<http://www.worldcat.org/oclc/54084501>>
<<http://elvis.slis.indiana.edu/irpub/DL/1997/pdf5.pdf>>
TopicModels/DDC: Mehler, Alexander; Waltinger, Ulli: Enhancing Document Modeling by Means of Open Topic Models: Crossing the Frontier of Classification Schemes in Digital Libraries by Example of the DDC. Library Hi Tech, Vol. 27, Issue 4, 2009, pp. 520 – 539.
- 9 Lewis, David D.: Evaluating Text Categorization. Proceedings of the Workshop on Speech and Natural Language. Association for Computational Linguistics. Morristown, NJ, USA, 1991, pp. 312 - 318.
<<http://portal.acm.org/citation.cfm?id=112471>>
- 10 Larson, Ray R.: Experiments in Automatic Library of Congress Classification. Journal of the American Society for Information Science, Vol. 43, pp. 130 - 148; zitiert nach (Wang 2009), S. 2281.
- 11 Zurück auf Los <<http://blog.bib.uni-mannheim.de/Classification/?p=30>>
- 12 Weitere Version mit umfangreicheren Erläuterungen und Quellenangaben: <<http://taipan.dyndns.org/~ul/dialog2010.pdf>>