

**Günter Bachelier**

Polyrepräsentation,  
Relevanz-Approximation und aktives  
Lernen im Vektorraummodell des  
Information-Retrievals

**Doktorarbeit / Dissertation**

## **Bibliografische Information der Deutschen Nationalbibliothek:**

Bibliografische Information der Deutschen Nationalbibliothek: Die Deutsche Bibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de/> abrufbar.

Dieses Werk sowie alle darin enthaltenen einzelnen Beiträge und Abbildungen sind urheberrechtlich geschützt. Jede Verwertung, die nicht ausdrücklich vom Urheberrechtsschutz zugelassen ist, bedarf der vorherigen Zustimmung des Verlanges. Das gilt insbesondere für Vervielfältigungen, Bearbeitungen, Übersetzungen, Mikroverfilmungen, Auswertungen durch Datenbanken und für die Einspeicherung und Verarbeitung in elektronische Systeme. Alle Rechte, auch die des auszugsweisen Nachdrucks, der fotomechanischen Wiedergabe (einschließlich Mikrokopie) sowie der Auswertung durch Datenbanken oder ähnliche Einrichtungen, vorbehalten.

Copyright © 2001 Diplomica Verlag GmbH  
ISBN: 9783832465636

**Günter Bachelier**

**Polyrepräsentation, Relevanz-Approximation und aktives Lernen im Vektorraummodell des Information-Retrievals**

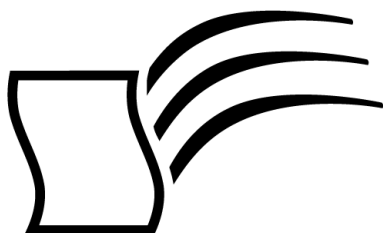


---

Günter Bachelier

# **Polyrepräsentation, Relevanz- Approximation und aktives Lernen im Vektorraummodell des Information-Retrievals**

**Dissertation / Doktorarbeit  
an der Universität des Saarlandes  
Fachbereich Informationswissenschaft  
Juni 2001 Abgabe**



***Diplom.de***

Diplomica GmbH \_\_\_\_\_  
Hermannstal 119k \_\_\_\_\_  
22119 Hamburg \_\_\_\_\_

Fon: 040 / 655 99 20 \_\_\_\_\_  
Fax: 040 / 655 99 222 \_\_\_\_\_

agentur@diplom.de \_\_\_\_\_  
www.diplom.de \_\_\_\_\_

ID 6563

Bachelier, Günter: Polyrepräsentation, Relevanz-Approximation und aktives Lernen im Vektorraummodell des Information-Retrievals

Hamburg: Diplomica GmbH, 2003

Zugl.: Saarbrücken, Universität, Dissertation / Doktorarbeit, 2001

---

Dieses Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere die der Übersetzung, des Nachdrucks, des Vortrags, der Entnahme von Abbildungen und Tabellen, der Funksendung, der Mikroverfilmung oder der Vervielfältigung auf anderen Wegen und der Speicherung in Datenverarbeitungsanlagen, bleiben, auch bei nur auszugsweiser Verwertung, vorbehalten. Eine Vervielfältigung dieses Werkes oder von Teilen dieses Werkes ist auch im Einzelfall nur in den Grenzen der gesetzlichen Bestimmungen des Urheberrechtsgesetzes der Bundesrepublik Deutschland in der jeweils geltenden Fassung zulässig. Sie ist grundsätzlich vergütungspflichtig. Zuwiderhandlungen unterliegen den Strafbestimmungen des Urheberrechtes.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

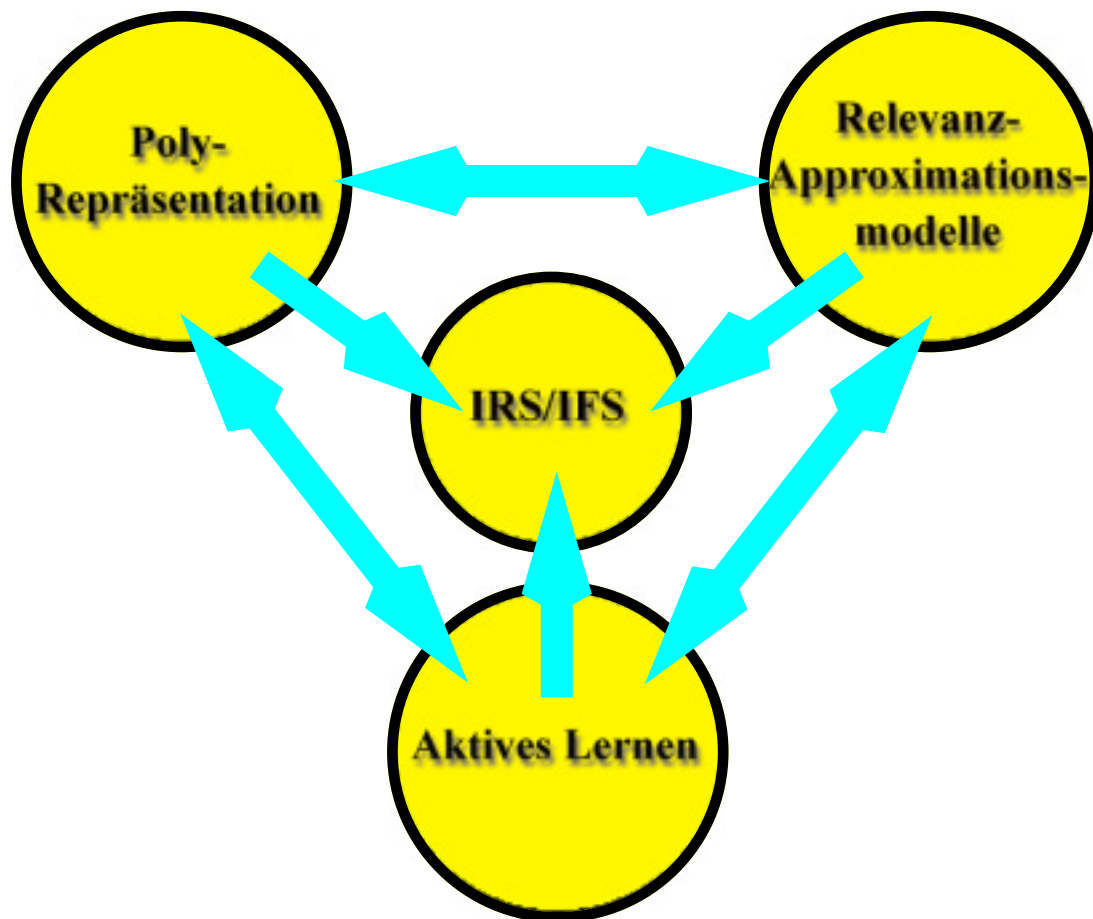
Die Informationen in diesem Werk wurden mit Sorgfalt erarbeitet. Dennoch können Fehler nicht vollständig ausgeschlossen werden, und die Diplomarbeiten Agentur, die Autoren oder Übersetzer übernehmen keine juristische Verantwortung oder irgendeine Haftung für evtl. verbliebene fehlerhafte Angaben und deren Folgen.

Diplomica GmbH

<http://www.diplom.de>, Hamburg 2003

Printed in Germany

# Polyrepräsentation, Relevanz-Approximation und aktives Lernen im Vektorraummodell des Information-Retrievals



von

Günter Bachelier

**Dissertation**

zur Erlangung des Grades eines Doktors der Philosophie der  
Philosophischen Fakultäten der Universität des Saarlandes

Juni 2001

Tag der Promotion: 31.01.2002

Dekan der Philosophischen Fakultät III: Herr Univ.-Prof. Dr. Ernst Löffler

Erstberichterstatter: Herr Univ.-Prof. Dr. Harald H. Zimmermann

Zweitberichterstatter: Herr Univ.-Prof. Dr. Werner Tack



## Vorwort

Konzipiert wurde diese Arbeit als elektronische Publikation in Form einer pdf-Datei mit einer Vielzahl interner und externer Links, sodass einige kurze Anwendungshinweise vorangestellt werden sollen.

Blaue, nicht unterstrichene Textbereiche repräsentieren Hyperlinks innerhalb der Arbeit, während blau unterstrichene Textbereiche auf externe Quellen (URLs) verweisen.

Eine größere Anzahl von pdf-Dateien, die im [Literaturverzeichnis](#) angegeben sind, können über die ACM Digital Library geladen werden. Der Zugang zur ACM Digital Library erfordert jedoch die Mitgliedschaft bei ACM, sowie eine Online-Anmeldung unter <http://www.acm.org/dl/>.

Einige Dateien, die im [Literaturverzeichnis](#) angegeben sind, besitzen den Typ XXX.djvu. Um entsprechende Dateien mit einem Browser zu betrachten, wird das Browser-Plugin "DjVuBrowser" benötigt, welches unter <http://www.djvuzone.org/> geladen werden kann (Stand Mai 2001).

Die URL-Quellen beziehen sich bis auf die Angaben des Typs XXX.djvu auf den Stand April/Mai 2001.

Danken möchte ich Prof. Dr. Harald H. Zimmermann, dass ich die dargestellten Themenbereiche frei auswählen und kreativ bearbeiten konnte, sowie David Cohn für die Diskussion bezüglich einiger meiner Ideen zum aktiven Lernen (Cohn (2000[75])).

Günter Bachelier M. A.

## Inhalt

1) Einleitung und Überblick .....	13
1.1) Information-Retrieval-Systeme als Spezialfall von Informationssystemen .....	13
1.2) Problemkomplexität des Information Retrievals .....	14
1.2.1) Hochdimensionale Zusammenhänge .....	15
1.2.2) Nicht-lineare und multimodale Zusammenhänge .....	15
1.2.3) Dynamische Zusammenhänge (nicht-stationäre Funktionen) .....	16
1.2.4) Unsicherheit (uncertainty) und Vagheit (fuzzyness) .....	18
1.2.5) Diversität der Agenten und ihre Ziele .....	19
1.2.6) Mehrziel-Anforderungen .....	20
1.3) Methodentransfer .....	21
1.4) Adaptive Informationssysteme .....	35
1.5) Einbettung externer Informationsbeschaffung in ein Modell der allgemeinen Intelligenz .....	37
1.6) Polyrepräsentation .....	43
1.6.1) Polyrepräsentation in IS und IRS .....	44
1.6.2) Inter- und Intraparadigmen-Polyrepräsentation .....	45
1.6.3) Gründe für das Vektorraummodell als Intraparadigmen-Polyrepräsentation .....	47
1.6.4) Gründe für Polyrepräsentation .....	48
1.6.4.1) Beschränkung endlicher Lernmengen .....	48
1.6.4.2) Fundamentale Beschränkung aller Repräsentationssprachen .....	49
1.6.4.3) Modellierung von Unsicherheit in Bezug zur Diversität der Agenten .....	51
1.7) Relevanz-Approximationsmodelle .....	52
1.8) Aktives Lernen .....	55
1.8.1) Passives und aktives Lernen .....	55
1.8.2) Geschlossene und offene Lernmenge .....	56
1.8.3) Direkte und indirekte Verfahren des aktiven Lernens .....	57
1.8.4) Effektivitäts- und Effizienz-Vergleich direkter und indirekter Verfahren .....	58
1.8.5) Relevanz- und Modell-Maximierungskriterium .....	58
1.8.6) Modell-Polyrepräsentation beim aktiven Lernen .....	60
1.9) Semantisches Netz der Beziehungen der verwendeten Ansätze .....	60
2) Methodische Grundlagen .....	62
2.1) Basis-Verfahren der stützpunkt-basierten Approximation .....	62
2.1.1) Vektorraum und Metrik .....	62
2.1.2) Approximation, Interpolation, Regression .....	63
2.1.2.1) Symbolische und stützpunktorientierte Approximation .....	66
2.1.2.2) Instanz- und prototypbasierte stützpunktorientierte Approximation .....	66
2.1.2.3) Framework des überwachten und unüberwachten Lernens .....	67
2.1.3) Local-Weighted-Regression .....	68
2.1.4) Sensorische-SOM (S-SOM) .....	70
2.1.5) Growing-Neural-Gas (GNG-SOM) .....	72

2.1.6)	Stimulus-Cluster-GNG-SOM (SC-GNG-SOM)	75
2.1.7)	Batch-Lernen in einer SC-GNG-SOM	78
2.1.8)	Aktivitätsausbreitung in GNG-Graphen	83
2.2)	Basis-Verfahren des Resamplings	90
2.2.1)	Stimulus-Bootstrap	90
2.2.2)	Restwert-Bootstrap	91
2.2.3)	Moving-Blocks-Bootstrap	91
2.3)	Modellqualität	94
2.3.1)	Unüberwachte SC-GNG-SOM-Qualität durch lokale Quantifizierungsfehler	94
2.3.2)	Überwachte SC-GNG-SOM-Qualität durch lokale MSE-Werte	95
2.3.3)	MSE-Integral und Bias-Varianz-Zerlegung	96
2.3.4)	Modellbewertung durch Varianz- und Bias-Integrale	97
2.3.5)	Output-, Bias- und Varianz-Approximationsmodelle	98
2.3.6)	Suche nach Inputvektoren mit extremalen Outputwerten	101
2.3.7)	Modellbewertung durch Momente höherer Ordnung	101
2.3.7.1)	Output- und Fehler-Moment bei einem Modell	102
2.3.7.2)	Output- und Fehler-Moment bei einer Modellmenge	103
2.3.7.3)	Schätzung von Output- und Fehler-Momenten	104
2.3.8)	Modellqualität durch Gewinnerlisten-Verfahren	107
2.4)	Hierarchische Strukturierung bei Mehr-Ziel-Optimierungen	111
2.4.1)	Pareto-kriterium und Paretomenge	112
2.4.2)	Pareto-Hierarchien	114
2.4.2.1)	Dominanz-Ranking	114
2.4.2.2)	Sukzessive Deaktivierung von Paretomengen	115
2.4.2.3)	Pareto-Wettkampf-Hierarchien	116
2.4.2.3.1)	Wettkampfoperation	117
2.4.2.3.2)	Wettkampf-Hierarchie	118
2.4.2.3.3)	Pareto-Wettkampf-Hierarchie als Spezialisierung einer Wettkampf-Hierarchie	119
2.4.3)	Abbruchkriterium bei Mehr-Ziel-Optimierung	120
2.5)	Verwendete Modelle der Evolutions-Strategien	126
2.5.1)	Ein-Ziel-Optimierung durch $\rho$ -geschlechtliche $(\mu, \lambda)$ - und $(\mu + \lambda)$ -ES	126
2.5.2)	Mehr-Ziel-ES	130
2.6)	Intervall-Selektions-Operatoren	135
2.6.1)	Rangfolge durch Ordnen nach einem ausgewählten Punkt im Intervall	136
2.6.2)	Selektion durch Zugehörigkeitsfunktionen	138
2.6.3)	Dominanzfunktion auf der Basis von Intervallen	139
3)	Mono- und Polyrepräsentation im vektorraumbasierten Information-Retrieval	142
3.1)	Information Retrieval Systeme	142
3.2)	Dokument als Zeichensequenz	146
3.2.1)	Dokument-Monorepräsentation	146
3.2.2)	Dokument-Polyrepräsentation	147

3.3) Merkmale als Zeichensequenz .....	151
3.3.1) Monorepräsentation von Merkmalen .....	151
3.3.2) Polyrepräsentation von Merkmalen .....	152
3.4) Indexierung .....	153
3.4.1) Monorepräsentation der Indexierung .....	153
3.4.2) Polyrepräsentation der Indexierung .....	155
3.5) Merkmalsgewichtungsmodelle .....	158
3.5.1) Grundlegende Merkmalsgewichtungsmodelle .....	158
3.5.2) Mono- und Polyrepräsentation der Indexierung im Kontext der Merkmalsgewichtungsmodelle .....	161
3.6) Retrieval .....	163
3.6.1) Query als Zeichensequenz .....	165
3.6.1.1) Query-Monorepräsentation .....	165
3.6.1.2) Query-Polyrepräsentation .....	165
3.6.1.2.1) Polyrepräsentation durch multiple Queryformulierung eines Agenten .....	165
3.6.1.2.2) Polyrepräsentation durch kollaborative Queryformulierung einer Agentengruppe .....	166
3.6.1.2.3) Polyrepräsentation durch Moving-Blocks-Bootstrap .....	166
3.6.1.2.4) Polyrepräsentation durch Mutations-Operationen .....	168
3.6.1.2.5) Polyrepräsentation durch Markov-Prozesse .....	168
3.6.1.2.6) Polyrepräsentation durch Rekombinations-Operationen .....	169
3.6.1.2.7) Polyrepräsentation durch GNG-SOM-Merkmalsgraphen .....	172
3.6.2) Query-Indexierung und Queryvektor-Mono- und Polyrepräsentation .....	173
3.6.2.1) Queryvektor-Monorepräsentation .....	173
3.6.2.2) Queryvektor-Polyrepräsentation .....	173
3.6.2.2.1) Query-Polyrepräsentation und Indexierungsfunktions- Monorepräsentation .....	173
3.6.2.2.2) Query-Monorepräsentation und Indexierungsfunktions- Polyrepräsentation .....	174
3.6.2.2.3) Query-Monorepräsentation und stochastische Indexierungsfunktion .....	174
3.6.2.2.4) Query-Polyrepräsentation und Queryvektor-Reproduktions- Operationen .....	175
3.6.3) Retrievalstrategien bei einer Dokumentvektor-Monorepräsentation .....	176
3.6.3.1) Dokumentvektor-Monorepräsentation und Queryvektor- Monorepräsentation .....	179
3.6.3.2) Dokumentvektor-Monorepräsentation und Queryvektor- Polyrepräsentation .....	183
3.6.4) Retrievalstrategien bei einer Dokumentvektor-Polyrepräsentation .....	187
3.6.5) Retrievalstrategien bei einer Retrievalregion-Mono- und -Polyrepräsentation .....	189

3.6.6) Retrievalstrategien mit positiven und negativen Queries und Queryvektoren .....	192
3.6.6.1) Monorepräsentation von positiven und negativen Queryvektoren .....	192
3.6.6.2) Polyrepräsentation von positiven und negativen Queryvektoren .....	195
3.7) Clusterung in IRS .....	198
3.7.1) Allgemeine Objekt- und Objektvektoren-Clusterung .....	199
3.7.2) Dokumentvektoren-Clusterung .....	201
3.7.3) Merkmalsvektoren-Clusterung .....	203
3.7.4) Integrierte Dokumentvektoren und Merkmalsvektoren-Clusterung .....	204
3.7.5) Cluster-Retrieval-Strategien .....	206
3.8) Indexierung und Retrieval mit GNG-SOM-Modellen am Beispiel unabhängiger Merkmals- und Dokument-Graphen .....	211
3.8.1) Aufbau unabhängiger Graphen .....	213
3.8.2) Einfache Cluster-Retrieval-Strategien mit Dokumentvektoren-Graph .....	215
3.8.3) Cluster-Retrieval-Strategien mit positiven und negativen Queryvektoren .....	218
3.8.4) Triangulation positiver und negativer Queryvektoren .....	221
3.8.5) Retrieval-Strategien mit Query-Modifikation .....	225
3.8.6) Queryvektor-Polyrepräsentation in Merkmalsgraphen .....	227
3.9) Relevanz-Feedback in IRS .....	230
3.9.1) Relevanzbegriff und Relevanzproblematik .....	232
3.9.1.1) Ähnlichkeits-Relevanz .....	232
3.9.1.2) Problemlösungs-Relevanz .....	232
3.9.1.3) Modellierung des Problemlösungsprozesses durch einen Zustandsraum .....	235
3.9.1.4) Reformulierungs-Relevanz .....	236
3.9.1.5) Irrelevant vs. irreführend .....	237
3.9.2) Queryvektor-Relevanz-Feedback .....	238
3.9.2.1) Queryvektor-Feedback bei unklassifizierten Dokumentvektoren .....	238
3.9.2.2) SOM-Adaption beim Queryvektor-Feedback .....	243
3.9.2.3) Stochastische Adaptions-Operationen beim Queryvektor-Feedback .....	245
3.9.2.4) Post-Retrieval-Operationen beim Queryvektor-Feedback .....	247
3.9.2.5) Queryvektor-Splitting .....	249
3.9.2.6) Queryvektor-Polyrepräsentation beim Queryvektor-Feedback .....	251
3.9.2.7) Queryvektor-Feedback bei Dokumentvektoren-GNG-SOMs .....	255
3.9.2.8) Queryvektor-Feedback mit positiven und negativen Queryvektoren .....	258
3.9.2.9) Queryvektoren-Trajektorie .....	262
3.9.3) Dokumentvektor-Relevanz-Feedback .....	266
3.9.3.1) Dokumentvektor-Feedback bei unklassifizierten Dokumentvektoren .....	266
3.9.3.2) Dokumentvektor-Feedback durch SOM-Adaption .....	270
3.9.3.3) Dokumentvektor-Feedback bei SC-GNG-SOMs .....	274
3.9.4) Gewichtsvektor-Relevanz-Feedback .....	277
3.9.5) Retrievalregion-Relevanz-Feedback .....	280

3.9.6) Indexierungsfunktion-Relevanz-Feedback .....	288
3.9.6.1) Detaillierte Beschreibung der Indexierungsfunktion .....	289
3.9.6.2) Indexierungsfunktionssuche nach dem Queryvektor-Feedback .....	291
3.9.6.2.1) Fitnessfunktion .....	291
3.9.6.2.2) Strategien zur Effizienzverbesserung .....	293
3.9.6.3) Indexierungsfunktionssuche parallel zum Queryvektor-Feedback .....	294
3.9.6.3.1) Queryvektor-Feedback ohne Reindexierung .....	295
3.9.6.3.2) Queryvektor-Feedback mit Reindexierung .....	297
3.9.7) Reformulierungs-Relevanz-Feedback .....	298
3.9.7.1) Reformulierung der Query .....	299
3.9.7.1.1) Relevanzwerte aus Queries bzw. Queryvektoren .....	300
3.9.7.1.2) Direkte Frage nach Reformulierungs-Relevanzwerten .....	304
3.9.7.1.3) Relevanzwerte aus Queryvektoren und Bewertung .....	305
3.9.7.2) Reformulierung anderer Texttypen wie der Problembeschreibung .....	308
3.9.8) Gleichzeitige Modifikation mehrerer Repräsentationen am Beispiel von Queryvektor- und Dokumentvektor-Feedback .....	309
 4) Relevanz-Approximation in Mono- und Polyrepräsentations-IRS .....	312
4.1) Approximationsmodelle mit reellen Relevanzwerten .....	312
4.1.1) Binäre und reelle Relevanzwerte .....	312
4.1.2) Ranking und Distanz-Relevanzfunktion .....	313
4.1.3) Relevanz-Klassifikations- und Approximationsmodelle .....	318
4.2) Feedback mit reellen Relevanzbewertungen bei unklassifizierten Dokumentvektoren .....	320
4.2.1) Queryvektor-Adaption bei reellen Relevanzbewertungen .....	320
4.2.1.1) Adaption bei Queryvektor-Monorepräsentation .....	321
4.2.1.2) Adaption bei Queryvektor-Polyrepräsentation .....	324
4.2.1.3) Adaption bei positiven und negativen Queryvektoren .....	325
4.2.2) Feedback mit Relevanz-Approximationsmodell ohne Veränderung des Queryvektors .....	325
4.2.3) Feedback mit Approximationsmodell und nachträglicher Adaption des Queryvektors .....	328
4.2.4) Effizienzsteigerung des Modells ohne Queryvektor-Veränderung .....	328
4.2.4.1) Effizienzsteigerung durch Distanz- bzw. Kernel-Matrix .....	328
4.2.4.2) Effizienzsteigerung durch Einschränkung der Grundmenge .....	329
4.2.4.2.1) Einschränkung durch $\epsilon$ -Umgebung .....	329
4.2.4.2.2) Einschränkung durch GNG-SOM-Repräsentation .....	331
4.2.5) Clusterung der Gesamtergebnismenge durch GNG-SOM .....	334
4.2.6) Prototypbasiertes GNG-SOM-Approximationsmodell aus Gesamtergebnismenge .....	337
4.2.7) Combining-Strategie bei Ergebnismengenbildung .....	340

4.3) Feedback mit reellen Relevanzbewertungen bei Dokumentvektoren-GNG-SOMs .....	340
4.3.1) Dokumentvektoren-GNG-SOM mit instanzbasiertem Modell .....	341
4.3.2) Dokumentvektoren-GNG-SOM mit prototypbasiertem Modell .....	344
4.3.3) Dokumentvektoren-GNG-SOM mit Nachadaption .....	347
4.3.3.1) Nachadaption ohne Wachstumsoperationen .....	348
4.3.3.2) Nachadaption mit Wachstumsoperationen .....	349
4.4) Relevanz-Approximationsmodell-Polyrepräsentation .....	352
4.4.1) Polyrepräsentation bei instanzbasierten Approximationsmodellen .....	353
4.4.1.1) Approximationsmodell-Polyrepräsentation durch Queryvektor-Polyrepräsentation .....	353
4.4.1.2) Approximationsmodell-Polyrepräsentation durch Bootstrap-Verfahren .....	356
4.4.2) Polyrepräsentation bei prototypbasierten Approximationsmodellen .....	357
4.4.2.1) Polyrepräsentierte Prototyp-Modelle bei unklassifizierten Dokumentvektoren .....	358
4.4.2.1.1) Unabhängiger Aufbau von Prototyp-Modellen durch Stimulus-Bootstrap .....	359
4.4.2.1.2) Unabhängiger Aufbau von Prototyp-Modellen durch Neuronen-Bootstrap .....	361
4.4.2.1.3) Abhängiger Aufbau von Prototyp-Modellen durch Stimulus-Bootstrap .....	362
4.4.2.1.3.1) Beibehaltung von Bootstrap-GNG-SOMs .....	363
4.4.2.1.3.2) Iterations-spezifische Neuableitung von Bootstrap-GNG-SOMs .....	366
4.4.2.1.3.3) Bootstrap-GNG-SOMs durch Aktualisierung von Relevanzschätzungen .....	368
4.4.2.1.3.4) Relevanzschätzungs-Dichtefunktion .....	371
4.4.2.2) Polyrepräsentierte Prototyp-Modelle bei klassifizierten Dokumentvektoren .....	374
4.4.2.2.1) Adaption der Stützpunkte im Relevanzraum und Erhaltung im DVR .....	377
4.4.2.2.2) Adaption der Stützpunkte im DVR und im Relevanzraum .....	382
4.4.2.2.3) Adaption mit Wachstum der Stützpunkte im DVR und Relevanzraum .....	385
4.5) Nutzung von Ergebnissen vergangener Interaktionen .....	388
4.5.1) Selektionsverfahren für Interaktionsobjekte bei Mono- und Polyrepräsentation .....	390
4.5.2) Nutzung von Stimulismengen vergangener Interaktionen .....	393
4.5.2.1) Ergebnismengen durch lokale Operationen in T .....	393
4.5.2.1.1) Monorepräsentation von Relevanzwerten .....	393
4.5.2.1.2) Polyrepräsentation von Relevanzwerten .....	395

4.5.2.2)	Ergebnismengen durch Übernahme aus ausgewählten Interaktionssmengen .....	396
4.5.2.2.1)	Monorepräsentation der Interaktionsmenge .....	397
4.5.2.2.2)	Polyrepräsentation der Interaktionsmenge .....	398
4.5.3)	Nutzung von nachadaptierten Queryvektoren vergangener Interaktionen .....	400
4.5.3.1)	Monorepräsentation der Nutzung nachadaptierter Queryvektoren .....	401
4.5.3.2)	Polyrepräsentation der Nutzung nachadaptierter Queryvektoren .....	402
4.5.4)	Nutzung von Relevanz-Approximationsmodellen vergangener Interaktionen .....	403
4.5.4.1)	Monorepräsentation der Approximationsmodell-Nutzung .....	404
4.5.4.2)	Polyrepräsentation der Approximationsmodell-Nutzung .....	405
4.5.5)	Nutzung von Suchregionen vergangener Interaktionen .....	407
4.6)	Korrektur der Relevanzschätzungen um Fehlerschätzungen .....	410
4.6.1)	Instanzbasierte Fehlermodelle .....	411
4.6.1.1)	Monorepräsentation von instanzbasierten Fehlermodellen .....	411
4.6.1.2)	Polyrepräsentation von instanzbasierten Fehlermodellen .....	413
4.6.2)	Prototypbasierte Fehlermodelle .....	414
4.6.2.1)	Monorepräsentation von prototypbasierten Fehlermodellen .....	414
4.6.2.2)	Polyrepräsentation von prototypbasierten Fehlermodellen .....	416
4.7)	Unterschiedliche Gewichtung von Relevanzmaximierung und Modellaufbau .....	418
5)	Aktives Lernen in Mono- und Polyrepräsentations-IRS .....	423
5.1)	Passives und aktives Lernen .....	424
5.1.1)	Passives Lernen .....	425
5.1.2)	Aktives Lernen bei einer geschlossenen Stimulusmenge .....	427
5.1.3)	Aktives Lernen bei einem Stimulusstrom .....	431
5.1.4)	Aktives Lernen bei einer offenen Stimulusmenge .....	435
5.2)	Indirekte und direkte Verfahren beim Modell-Maximierungskriterium .....	439
5.2.1)	Indirekte Verfahren .....	440
5.2.1.1)	Selektionskriterien bei indirekten Verfahren .....	440
5.2.1.2)	Allgemeine Selektion durch fehlende Übereinstimmung .....	441
5.2.1.3)	Outputvarianz-Maximierung .....	442
5.2.1.4)	Bias- und Varianz-Maximierung bei klassifizierten Dokumentvektoren .....	443
5.2.1.4.1)	Unabhängige Listenbildung .....	448
5.2.1.4.2)	Abhängige Listenbildung .....	450
5.2.2)	Direkte Verfahren am Beispiel Optimal-Experiment Design .....	452
5.2.2.1)	Bias-Quadrat-Integral-Minimierung .....	452
5.2.2.2)	Output-Varianz-Integral-Minimierung .....	458
5.2.2.3)	Kombinierte Bias-Quadrat- und Output-Varianz-Integral-Minimierung .....	460



5.2.3) Effizienzverbesserungen bei direkten Verfahren .....	464
5.2.3.1) Eigenschafts-Integrale mit weniger Stützpunkten durch Fehlerauswahl .....	464
5.2.3.2) Eigenschafts-Integrale mit weniger Stützpunkten durch Häufigkeitsverteilung .....	467
5.2.3.3) Deterministische Integration im average case setting .....	469
5.2.3.4) Weniger Eigenschafts-Integrale durch Kandidatencluster und Approximation .....	470
5.2.3.5) Kombination von stetiger und diskreter Vorgehensweise .....	475
5.2.3.6) Neurone als Stützpunkte der Approximation und der Integration .....	476
5.2.4) Effektivitätsverbesserungen bei direkten Verfahren durch zentrale Momente .....	484
5.3) Integration von Relevanz- und Modell-Maximierungskriterium .....	486
5.3.1) Erzeugung einer gemeinsamen tertiären Ergebnisliste .....	486
5.3.1.1) Dokumentvektoren als Komponenten in beiden Listen .....	488
5.3.1.2) Dokumentvektormengen als Komponenten in beiden Listen .....	490
5.3.1.3) Dokumentvektoren als Komponenten der ersten und Dokumentvektormengen als Komponenten der zweiten Liste .....	491
5.3.2) Erzeugung zweier tertiären Ergebnislisten .....	493
5.3.3) Erzeugung dreier tertiären Ergebnislisten .....	494
5.4) Lösungsansatz des Kombinatorikproblems bei direkten Verfahren durch die Integration eines Output- und eines Modell-Maximierungskriteriums .....	495
5.4.1) Output-Maximierung .....	495
5.4.2) Vorstrukturierung der Kandidatenmenge .....	496
5.4.3) Bildung von Kandidatenteilmengen .....	496
5.4.4) Direktes aktives Lernen bei vorstrukturierten Kandidatenteilmengen .....	497
 6) Zusammenfassung .....	 501
 Verzeichnis ausgewählter Symbole .....	 510
Abbildungsverzeichnis .....	520
Literaturverzeichnis .....	524

## 1) Einleitung und Überblick

## 1.1) Information-Retrieval-Systeme als Spezialfall von Informationssystemen

Informationssysteme (IS) können durch einen 5- bzw. 7-Tupel beschrieben werden (Panyr (1986a: 22[247])):

$$IS = (A, W, Q, I, E). \text{ bzw. } IS = (A, W, Q, I, E, U, D). \quad (1)$$

- A: Inputfunktion (Erschließungsfunktion, Lernfunktion) zum Aufbau der internen Repräsentationen.
- W: Interne Repräsentationen (Wissensbasis).
- Q: Inputmenge als Menge aller zugelassenen Inputkonfigurationen (Problemformulierung, Suchfrage)
- I: Outputfunktion (Inferenzfunktion, Retrievalfunktion und Relevanz-Feedbackfunktion).
- E: Outputmenge als Menge aller möglichen Outputkonfigurationen (Problemlösung, Systemvorschlag).
- U: Updatefunktion der internen Repräsentationen.
- D: Dialogkomponente.

Information-Retrieval-Systeme (IRS) können als Spezialfall eines IS beschrieben werden, indem die einzelnen Komponenten des Tupels (A, W, Q, I, E) spezifiziert werden, was im Kontext des Standard-Retrieval-Prozesses in einem Vektorraummodell geschehen soll. Gegeben ist zu einem Zeitpunkt  $t$  eine Dokumentmenge  $D^t$ , die durch eine Dokument-Indexierungsfunktion  $A_{IR(D)}$  auf eine Dokumentvektorenmenge  $DVM^t$  abgebildet wurde. Die einzelnen Dokumentvektoren  $x_i$  sind Element eines metrischen,  $n^t$ -dimensionalen Dokumentvektorraumes  $DVR$ , mit  $n^t$  als der Anzahl der Merkmale (Features), auf der die Indexierung basiert und die in der Menge  $F^t$  zusammengefasst werden. Der Dokumentvektorraum wird allgemein als Teilraum von  $R^{n(t)}$  beschrieben, z.B. durch  $[0, 1]^{n(t)}$ . Die Query-Indexierungsfunktion  $A_{IR(Q)}$  wird vereinfachend definiert als Abbildung aus der Menge  $Q(\Theta)$  der möglichen bzw. zugelassenen Queries über einem endlichen Alphabet  $\Theta$ , in  $DVR$ . Es folgt die Anwendung der Retrieval-Funktion, die abhängig ist von der momentanen Dokumentvektorenmenge  $DVM^t$ , dem Queryvektor  $q_i^t$  und dem metrischen Dokumentvektorraum  $DVR$  mit seinen definierenden Eigenschaften, wobei hier ausschließlich die Metrik  $d_{DVR}$  betrachtet wird. Sei  $\Gamma_{DVR}$  die Menge aller Metriken, die in einem Dokumentvektorraum  $DVR$  angewendet werden können, ohne dass hier auf die Definition der Metrik eingegangen werden soll (siehe Abschnitt 3.6.3). Die Retrieval-Funktion kann somit spezifiziert werden als eine Abbildung der Potenzmenge  $P^{DVM(t)}$  der Dokumentvektorenmenge  $DVM^t$ , dem  $DVR$  und  $\Gamma_{DVR}$  auf  $P^{DVM(t)}$ , indem das Tripel aus  $DVM^t$ , dem Queryvektor  $q_i^t$  und eine Metrik  $d_{DVR}$  auf die query-abhängige Ergebnis-Dokumentvektorenmenge  $DVM_i^t$  abgebildet wird. D.h. die Retrieval-Funktion besitzt die allgemeine Form  $ret(DVM^t, q_i^t, d_{DVR})$  bzw.  $ret(DVM^t, q_i^t, d_{DVR}, \epsilon)$ , wenn eine einfache Best-Match-Retrievalstrategie betrachtet wird, bei der alle Dokumentvektoren aus  $DVM^t$  selektiert werden, deren Abstand von  $q_i^t$  kleiner-gleich einer Distanzschwelle  $\epsilon \in R^+$  ist. Der letzte Schritt besteht in der Erzeugung der Dokumentmenge  $D_i^t$ , die zu der Ergebnismenge  $DVM_i^t$  korrespondiert. Vereinfachend wurde auf die Beschreibung einer Ranking-Funktion verzichtet, die aus  $DVM_i^t$  eine geordnete Liste von Dokumentvektoren erzeugt. Wird diese Darstellung des Retrievalprozesses als Grundlage verwendet, so ergibt sich die grundlegende Beschreibung eines IRS ohne eine Relevanz-Feedbackfunktion als Element von I durch das folgende Tupel:

$$IRS = ((A_{IR(D)}, A_{IR(Q)}), (D^t, DVM^t, F^t, N_{DV}^t), Q(\Theta), ret(DVM^t, q_i^t, d_{DVR}), \{P^{DVM(t)}, P^{D(t)}\}). \quad (2)$$

- $A = (A_{IR(D)}, A_{IR(Q)})$ : Dokument-Indexierungsfunktion, Query-Indexierungsfunktion.
- $W = (D^t, DVM^t, F^t, (N_{DV}^t))$ : Dokumentmenge, Dokumentvektorenmenge, Merkmalsmenge, (Klassifikation  $N_{DV}^t$ ), (Wörterbücher, Thesauri u.ä. bleiben hier unberücksichtigt)
- $Q = Q(\Theta)$ : Menge aller Zeichensequenzen über einem endlichen Alphabet als Menge aller möglichen bzw. zugelassenen Queries.
- $I = \text{ret}(DVM^t, q_i, d_{DVR})$ : Retrievalfunktion als Funktion der Dokumentvektorenmenge, eines Queryvektors und der Metrik des Dokumentvektorraums DVR.
- $E \in \{P^{DVM^t}, P^{D(t)}\}$ : Menge aller möglichen System-Outputs als Potenzmenge der Dokumentvektorenmenge oder Potenzmenge der Dokumentmenge.

## 1.2) Problemkomplexität des Information Retrievals

Das Ziel des Information Retrievals (IR, (Rijsbergen (1979[279]), Salton & McGill (1987[294]), Grossman & Frieder (1999[154])) ist die Präsentation einer Teilmenge oder geordneten Liste (Ranking) von Dokumenten aus einer Gesamt-Dokumentmenge als Output auf einen Input in Form einer Anfrage (Query) durch einen Nutzer (Agenten). Es wird angenommen, dass die Anfrage ein spezielles Informationsbedürfnis des Agenten repräsentiert, sodass die Ergebnis-Dokumentmenge zur Befriedigung dieses Informationsbedürfnisses beitragen soll. Die internen Selektionsoperationen des Information Retrieval Systems (IRS), welche eine Ergebnis-Dokumentmenge erzeugen, sollen derart sein, dass nach der Integration des Inhaltes der nachgewiesenen Dokumente in das kognitive System des Agenten, das Informationsbedürfnis gegenüber dem Zustand vor der Integration kleiner wird bzw. sich niveliert. Eine alternative Sicht geht nicht von der Verringerung des Informationsbedürfnisses aus, sondern allgemein von dessen Transformation in ein neues Informationsbedürfnis, d.h. der Constraint eines monoton fallenden Maßes an Informationsbedürfnis wird aufgegeben. Dadurch wird die Situation modellierbar, dass ein tiefer gehendes Verständnis eines Problemes andere oder sogar mehr Fragen aufwirft, sodass ein anderes bzw. größeres Informationsbedürfnis aus der Wissenserweiterung des Agenten durch Integration von Informationen aus nachgewiesenen Dokumenten folgen kann (siehe auch Swanson (1987[328])).

Das Ziel der IR-Forschung ist die Suche nach adäquaten und implementierbaren Modellen des IR-Prozesses, wobei Adäquatheit durch unterschiedliche Performancemaße wie z.B. Recall und Precision operationalisiert wird. In den letzten 40 Jahren IR-Forschung hat sich das Problem der Performanceverbesserung bezüglich der verwendeten Maßen als ausgesprochen hartnäckig dargestellt, was sich im Kontext des IR im Internet noch verschärft hat. Es stellt sich die grundsätzliche Frage, warum das IR ein so komplexes Problem ist, wofür die Kombination einer Reihe von Eigenschaften verantwortlich ist, die im Nachfolgenden skizziert werden sollen:

- 1) Hochdimensionale Zusammenhänge.
- 2) Nicht-lineare und multimodale Zusammenhänge.
- 3) Dynamische Zusammenhänge (nicht-stationäre Funktionen)
- 4) Quellen von Unsicherheit (uncertainty) und Vagheit (fuzzyness).
- 5) Diversität der Agenten und ihrer Ziele
- 6) Mehrziel-Anforderungen.

### 1.2.1) Hochdimensionale Zusammenhänge

Natürlichsprachliche Texte und insbesondere Fachtexte zeichnen sich durch eine große Anzahl darin enthaltener Terme aus, sodass zur Repräsentation dieser Texte ein entsprechend hochdimensionaler Parameterraum (Dokumentvektorenraum) notwendig ist. Dies muss kombiniert betrachtet werden mit einer großen Anzahl von Dokumenten, was im Fall vektorraum-basierter IR in einer hochdimensionalen Term-Dokument-Matrix mündet.

Der Versuch, durch unterschiedliche Verfahren eine Dimensionsreduktion dieser Räume zu erreichen, wurde dementsprechend in der IR-Forschung früh unternommen, wobei die Clusterung von Dokumenten und Termen im Zentrum der Bemühungen stand (Sparck Jones (1971[319]), Crouch (1972[84]), Robertson (1977a[283], b[284], c[285], 1979[286]), Salton (1971[293]), Panyr (1986a[247]), Bachelier (1995[14])). Neuere Verfahren versuchen mit einer Hauptkomponenten-Analyse (Principal Component Analysis (PCA); Oja (1993[240]), Diamantaras & Kung (1996[94]), Lee (1998[200])) eine Dimensionsreduktion zu erreichen wie das Latent Semantic Indexing (LSI; Dumais et al. (1988[100]), Furnas et al. (1988[136]), Bartell et al. (1992[27]), Soboroff et al. (1998[315]), Hofmann (1999[165]), Ando (2000[11])). Hierzu zählen auch konnektionistische Verfahren, die nicht-lineare Faktorenanalysen durch Self-Organizing-Maps (SOMs; Ritter et al. (1991[282]), Kohonen (1989[184], 1995[185]), Bachelier (1998a[15])) approximieren und dies auf den Fall des IR abstimmen (Scholtes (1993[301]), Bachelier (1995[14])).

In der Literatur zum Machine Learning und zur Optimierung ist die Problematik hochdimensionaler Zusammenhänge unter dem Begriff curse-of-dimensionality bekannt, der von Bellman (1967[41]) eingeführt wurde (siehe z.B. Venkatesh et al. (1992[348]), Warwick & Kárny (1997[357])). Bezeichnet wird damit das exponentielle Wachstum von Regionen eines Suchraumes bei einem linearen Wachstum der Dimensionsanzahl, womit der Suchaufwand in dem entsprechenden Raum bei deterministischen Verfahren ebenfalls exponentiell steigt. Als Lösungsansatz gegen den curse-of-dimensionality können Verfahren zur Randomization (siehe z.B. Traub et al. (1988[338]), Wozniakowski (1996a[372], b[373])) eingesetzt werden, wie z.B. Monte Carlo Verfahren (Fishman (1995[114])), die unabhängig von der Dimensionalität des zugrunde liegenden Raumes arbeiten. Erkauft wird diese Eigenschaft damit, dass die gefundenen Lösungen nur mit einer bestimmten Wahrscheinlichkeit in einem definierten Lösungsintervall liegen, d.h. dass die Lösung nur für den sogenannten average-case gelten (siehe auch Abschnitt 5.2.3)).

### 1.2.2) Nicht-lineare und multimodale Zusammenhänge

Indexierungs- und Retrievalfunktionen in vektorraumbasierten Modellen verwenden nicht-lineare, reelle Gewichtungen von Merkmalen und Dokumenten. Dies ergibt sich daraus, dass es sinnvoll ist, Gewichtungen auf ein reelles Intervall wie  $[0, 1]$  zu normieren, da die Interpretierbarkeit beliebig großer oder kleiner Gewichte verloren geht. Das gleiche gilt für Schätzungen der Relevanz durch ein IRS, (siehe Kapitel 4)) wobei reelle Relevanzschätzungen als nicht-lineare, multimodale Funktion der Distanzen von Dokumentvektoren ein Schwerpunkt der vorliegenden Arbeit bildet.

Unimodale Funktionen sind Funktionen mit genau einem Extremwert. Eine unimodale Relevanzfunktion würde bedeuten, dass an genau einem Punkt im Dokumentvektorenraum ein Relevanzmaximum existiert, wobei das Dokument mit dem am nächsten liegenden Dokumentvektor das relevanteste Dokument ist. Notwendig muss von einer monoton fallenden Relevanzfunktion um das Maximum einer unimodalen Funktion ausgegangen werden, d.h. je weiter ein Dokumentvektor vom Relevanzmaximum entfernt liegt, desto geringer wird dessen Relevanz bewertet. Die Modellierung des Retrievals mit einem Queryvektor und einer unimodalen Funktion ist die Standardsituation in den vektorraumbasierten Modellen, d.h. das Relevanzmaximum wird an der Stelle des Queryvektors bzw. des am nächsten liegenden Dokumentvektors angenommen, und allen anderen Dokumentvektoren wird explizit oder implizit ein Relevanzwert als monoton fallende Funktion der Distanz zu dem Maximum zugeordnet. Diese Vorgehensweise ist stark simplifizierend, da mit einer unimodalen Funktion alle Fälle nicht modellierbar werden, bei denen mehrere Cluster relevanter Dokumentvektoren existieren, zwischen denen nicht-relevante Dokumentvektoren liegen. Diese Fälle lassen sich nur mit multimodalen Funktionen modellieren, bei denen jeder Cluster relevanter Dokumentvektoren durch ein eigenes Maximum der Relevanzfunktion beschrieben werden kann. Diese Argumentation wird in Kapitel 4) verwendet, um die effektive Verwendung multimodaler, nicht-linearer Relevanz-Approximationsmodelle zu begründen.

Demgegenüber verwenden probabilistische Modelle nicht-lineare Verteilungsannahmen, welche die Indexierung und das Retrieval als nicht-lineare Prozesse kennzeichnet, jedoch werden nur unimodale Verteilungen verwendet, die nur durch eine Kombination im Sinne eines „Mixture Models“ (McLachlan & Basford (1987[214])) zu einer multimodalen Funktion aggregiert werden können.

Nicht-lineare und multimodale Zusammenhänge in hochdimensionalen Räumen zu modellieren, ist eine mathematisch anspruchsvolle Aufgabe und erfordert große Rechen- und Speicherressourcen, sodass effektive Verfahren im IR entsprechend aufwendig werden.

### 1.2.3) Dynamische Zusammenhänge (nicht-stationäre Funktionen)

Die offensichtliche Dynamik in IRS besteht in der Aufnahme neuer Dokumente und deren Integration in die gegebene Repräsentationsstruktur, die sich dadurch lokal bzw. global verändern kann. Insbesondere beim Vorliegen von Verfahren zur Dimensionsreduktion kann die Integration neuer Informationsobjekte gravierende Folgen besitzen, wobei eine dynamische Anpassung einer Clusterstruktur eine vergleichsweise problemlose Reorganisation darstellt (Crouch (1972[84])). Verändert sich die Anzahl der Objekte (Informationsobjekte oder Cluster) in einem Raum eines vektorraum-basierten IRS, so verändert sich die Dimension der Merkmals-Dokument-Matrix, und somit die Dimension des komplementären Raumes über sequentiell und rekursiv abhängige Merkmals- und Dokument-Graphen (siehe Abschnitt 3.8); siehe auch Bachelier (1995[14])). Dies erfordert globale Reorganisations-Operationen oder im Extremfall einen Neuaufbau der Strukturen in dem betreffenden Raum, sodass die Anzahl der dynamischen Anpassungen unter dem Constraint der gewünschten Aktualität minimiert werden sollte.

Durch die Aufnahme neuer Dokumente besteht immer die Chance, dass darin Terme enthalten sind, die bislang dem IRS unbekannt waren. Eine Aufnahme neuer Terme verändert die Dimension des Dokumen-

traumes, mit der Notwendigkeit der Reindexierung aller vorhandenen Dokumente und der Anpassung etwaig vorhandener Strukturierungen wie einer Klassifikation.

Diese dynamischen Prozesse der Reorganisation von Repräsentationsstrukturen machen die Retrieval-Funktion zu einer nicht-stationären Funktion, d.h. der gleiche Input kann zu unterschiedlichen Zeitpunkten einen unterschiedlichen Output erzeugen. Die nicht-stationäre Retrieval-Funktion ist von Ereignissen abhängig, die bezogen auf das IRS als extern zu bezeichnen sind, da das IRS keinen Einfluss auf die Produktion von Dokumenten und den darin enthaltenen Termen besitzt. Der einzige Einfluss besteht in der Entscheidung, ob neue Dokumente und Terme repräsentiert werden sollen. Ein IRS, das diese Entscheidung selbst treffen könnte, wäre autonom im Sinne autonomer Agenten, doch diese Entscheidung wird dem IRS ebenfalls von externen Instanzen in Form der Systembetreiber vorgegeben.

Es existiert ein Trade-of zwischen hochdimensionalen Zusammenhängen und dynamischen Funktionen beim vektorraum-basierten IRS in Bezug auf die Merkmalsanzahl. Werden keine Terme als Merkmale verwendet, sondern Zeichen-n-Gramme (Scholtes (1993[301]), Bachelier (1995[14])), so besitzt der Dokumentvektorenraum a priori die maximale Dimension, die somit nicht mehr verändert werden muss. Sei  $\Theta$  ein endliches Alphabet mit  $|\Theta|$  Zeichen (siehe Kapitel 3), so besitzt ein Dokumentvektorenraum auf der Basis von Zeichen-n-Grammen eine Dimensionsanzahl von  $|\Theta|^n$ . D.h. jede Dimension des Dokumentvektorenraumes repräsentiert ein mögliches Zeichen-n-Gramm, wobei der Wert, der einer Komponente eines Dokumentvektors zugeordnet wird, d.h. die Merkmalsgewichtung, eine nicht-lineare Funktion der Häufigkeit des Auftretens des entsprechenden n-Gramms in dem zugrunde liegenden Dokument ist. Die Stabilität der Dimension des Dokumentvektorenraums wird somit erkaufte durch seine große Anzahl an Dimensionen, wobei das Gesamtsystem dennoch dynamisch bleibt, da Dokumente neu indiziert werden, was die Verteilungsstruktur im Dokumentvektorenraum, sowie die Dimension des Merkmalsraumes verändert.

Allgemein wird der Umgang, die Modellierung und Optimierung und Control von nicht-stationären Funktionen als schwierig bewertet, wobei robuste, adaptive Verfahren wie evolutionäre Algorithmen notwendig werden (Goldberg & Smith (1987[145]), Grefenstette (1992[153]), Dasgupta & McGregor (1992[89]), Cobb & Grefenstette (1993[69]), Vavak et al. (1996[346]), Mori et al. (1997[224]), Mori et al. (1998[225]), Stagge (1998[322]), Dozier (2000[97]), Santo & Kita (2000[295])). In Kombination mit der Nicht-Linearität dieser Funktionen und den hochdimensionalen Räumen trägt diese Systemeigenschaft entscheidend zu der Gesamtkomplexität bei.

Die dynamische Aufnahme neuer Informationsobjekte koppelt ein IRS an die offene Welt (Popper (1994[262], 1995[263]), Mühlenbein (1994[229], 1995[230])), was explizite Folgen für die Methoden hat, da wahrscheinlichkeitstheoretische Induktionsmethoden nur in geschlossenen Welten gelten (Popper (1994: 376ff[262])). Ein IRS als offenes System (dissipatives System; Nicolis & Prigogine (1977[235])) in einer offenen Welt zu betrachten, wurde jedoch in der IR-Forschung bislang nicht als Modellierungsansatz verwendet.

### 1.2.4) Unsicherheit (uncertainty) und Vagheit (fuzzyness)

In Motro & Smets (1997[228]) werden Ursachen und Vorschläge zum Umgang mit Unsicherheit und Vagheit in Informationssystemen allgemein und IRS im speziellen (Turtle & Croft (1997[343])) gemacht (siehe auch Natke & Ben-Haim (1997[232])). Als Ursachen werden von Turtle & Croft (1997:191f[343]) folgende Faktoren genannt:

- a) Dokumentrepräsentation
- b) Repräsentation der Informationsbedürfnisse
- c) Retrievalfunktion
- d) Repräsentation im Kontext des Relevanz-Feedbacks

Die Indexierungsfunktion als Vorgang der Dokumentrepräsentation in vektorraumbasierten IRS bildet ein Dokument bzw. ein Merkmalshäufigkeitsvektor auf einen Dokumentvektor als Element eines metrischen Vektorraumes ab. Die Auswahl der Merkmale definiert die Struktur der Repräsentation, sodass Unsicherheiten bezüglich einer adäquaten Auswahl bestehen. Werden aus einer Dokumentmenge Terme automatisch extrahiert, so muss spezifiziert werden, welche Zeichensequenzen als Terme verwendet werden sollen, wodurch eine andere Unsicherheit bezüglich der Adäquatheit besteht.

Neben der Merkmalsauswahl ist die Merkmalsgewichtung die wesentliche Quelle der Unsicherheit, da eine Vielzahl von Vorschlägen für Gewichtungsfunktionen gemacht wurden (siehe z.B. Singhal (1997[312]), Zobel & Moffat (1998[376])). Welche Indexierungsfunktion für eine gegebene Dokumentmenge in Abhängigkeit von bestimmten Eigenschaften, die den Agenten als zukünftigen Nutzern zugeschrieben werden, hinreichend gut ist, kann nicht vorhergesagt werden (Zobel & Moffat (1998[376])).

Die Formulierung des Informationsbedürfnisses des Agenten in Form einer Query, führt weitere Unsicherheiten, Vagheiten und subjektive Faktoren wie der aktive Gebrauch eines Vokabulars ein und wird in der IR-Forschung am häufigsten zitiert. Wird das IR im Kontext eines Problemlösungsprozesses betrachtet (siehe Abschnitt 3.9.1)), und wird eine neuro-kognitive Sichtweise vertreten, die von funktionalen Hirnarealen ausgeht, so bedeutet die Queryformulierung eine Projektion interner Repräsentationsstrukturen aus einem funktionalen Problemlösungsareal auf Strukturen eines Sprachsyntheseareals. Es kann dabei nicht einmal von einer Abbildung als linkstotale, rechtseindeutige Relation ausgegangen werden, d.h. dass alle momentan existierenden Repräsentationsstrukturen aus dem Problemlösungsareal auf eine Struktur im Sprachsyntheseareal abgebildet werden können. Es besteht somit die Möglichkeit, dass Repräsentationsstrukturen nicht abgebildet werden und dass abgebildete Repräsentationen nicht umkehrbar abgebildet werden können. Sind z.B. bestimmte Vokabularelemente dem Agenten unbekannt, so können kognitive Strukturen darauf nicht abgebildet werden, was eine interne Quelle der Vagheit darstellt.

Wird ein interaktives Retrieval wie ein Relevanz-Feedback betrachtet (siehe Abschnitt 3.9), siehe auch Panyr (1987b[252]), Harman (1992[157]), Buckley & Salton (1995[63]), Cool et al. (1996[77]), Dunlop (1997[101]), Lundquist et al. (1997[203])), d.h. werden mehrere Dokumentlisten durch das IRS nachgewiesen, die der Agent bewerten soll, so kommen subjektive Faktoren, bezogen auf das Verstehen der Dokumente hinzu. D.h. die Abbildung von internen Strukturen aus dem Sprachanalyseareal auf Strukturen des Problemlösungsareals bildet eine weitere Quelle von Unsicherheiten und Vagheiten.

Weiterhin berücksichtigt werden muss, dass die Relationen bzw. Abbildungen zeitlich variabel sind, d.h. in allen Arealen verändern sich sowohl die Strukturen der Repräsentationen als auch die Anzahl der Repräsentationen. Entsteht in einem Areal eine neue Repräsentation durch interne Prozesse, so muss es in einer anderen Region, in die Repräsentationen abgebildet werden, keine oder noch keine korrespondierenden Strukturen geben, auf die abgebildet werden könnte. Denkbar wäre ein Prozess, dass zunächst auf nicht korrespondierende oder Default-Strukturen abgebildet wird, und dass im Rahmen einer Konsistenzherstellung neue Strukturen gebildet werden, auf die künftig eine Abbildung der entsprechenden Repräsentation erfolgt.

Als dritter Punkt wird von Turtle & Croft (1997[343]) die Retrievalfunktion angeführt, die als  $\text{ret}(\text{DVM}^t, q_i^t, d_{\text{DVR}}, \epsilon)$  definiert wurde als Funktion der Dokumentvektorenmenge, eines Queryvektors, der Metrik des DVRs und einem Distanzschwellenwert. Unsicherheit besteht bezogen auf die externe Festlegung der Indexierungsfunktion  $A_{\text{IR}(Q)}$  der Query, der Metrik und des Schwellenwertes. Die Metrik ist zudem zentraler Punkt jeder Rankingfunktion, die eine Ergebnis-Dokumentvektorenmenge auf eine geordnete Dokumentvektorenliste abbildet.

Als vierter Punkt werden von Turtle & Croft (1997[343]) Repräsentationen im Kontext des Relevanz-Feedbacks erwähnt. Relevanz-Feedback-Verfahren machen ein IRS zu einem adaptiven IRS, wobei interne Repräsentationen des IRS kurzfristig bzw. langfristig modifiziert werden. Wie diese Repräsentationen modifiziert werden, wird durch Adaptions-Funktionen definiert, die extern festgelegt werden müssen, und somit die Quelle von Unsicherheiten bezüglich der adäquaten Modifikation darstellen. Z.B. gehen Queryvektor-Relevanz-Feedback-Verfahren davon aus, dass die Query-Indexierungsfunktion  $A_{\text{IR}(Q)}$  nicht adäquat ist, sodass die Ursprungsquery mit Hilfe der in der ersten Iteration nachgewiesenen Dokumentvektoren und der Relevanzbewertungen durch den Agenten in einen neuen Queryvektor mit Hilfe einer Queryvektor-Adaptions-Funktion überführt wird. Die Verkettung von  $A_{\text{IR}(Q)}$  und der Adaptionsfunktion wird implizit als neue und adäquatere Query-Indexierungsfunktion verwendet. Neben dem Queryvektor-Relevanz-Feedback (siehe Abschnitt 3.9.2)) wird die Veränderung der Dokumentvektoren (Dokumentvektor-Relevanz-Feedback (siehe Abschnitt 3.9.3)) verwendet, während die Modifikation anderer Repräsentationsformen in der IR-Literatur unüblich ist (siehe Abschnitte 3.9.4) bis 3.9.8)).

### 1.2.5) Diversität der Agenten und ihre Ziele

Im Abschnitt über dynamische Zusammenhänge wurde bereits die Integration neuer Dokumente und Merkmale als Form der Kopplung eines IRS mit der offenen Welt beschrieben. Da das IRS keinen Einfluss auf die Produktion von Dokumenten besitzt, ist dies ebenfalls unter dem Gesichtspunkt der Unsicherheit interpretierbar, d.h. die Entwicklung der Repräsentationen in einem IRS sind durch diesen externen Faktor mit Unsicherheiten verbunden. Als weiterer externer Faktor wird die Diversität der Agenten betrachtet, d.h. es ist unsicher, welche Agenten mit welchen kognitiven Strukturen mit dem IRS zukünftig interagieren werden und welche Ziele sie damit verfolgen bzw. welche Probleme sie damit lösen wollen, wenn das IR im Kontext eines Problemlösungsprozesses betrachtet wird. Es werden immer neue Themen, Interessen, Ziele und Probleme auf der Ebene von Individuen und Gruppen gebildet, sowie immer neue Nutzergruppen aus unterschiedlichen Kulturen mit den IRS interagieren.



Weiterhin entstehen immer neue Mikro-Kulturen (Tribes) in denen Individuen mit gleichen Interessen, Zielen und Problemen sich global zusammenschließen, wobei ein Individuum zu einer Vielzahl unterschiedlicher Gruppen gehören kann. Ein solches Individuum kann nicht mehr mit einem Nutzermodell beschrieben werden, sondern es müsste kontextabhängig, d.h. in Abhängigkeit von der Gruppe, zu der das Individuum sich zum Zeitpunkt der Interaktion zugehörig fühlt, Nutzermodelle erzeugt werden, was eine Form von Polymorphismus darstellt. Durch diese Veränderungen wird die Diversität der Nutzer nicht nur in den Kontext der Unsicherheit, sondern auch in den Kontext dynamischer Zusammenhänge, d.h. nicht-stationärer Funktionen gestellt. Eine Polyrepräsentation von Nutzermodellen wird somit ein immer wichtig werdender Faktor, zumal eine wachsende Geschwindigkeit, in der sich die Diversität der Themen, Interessen, Ziele und Probleme vergrößert, beobachtbar ist.

Weiterhin kann ein Nutzer in einem engen Zeitintervall verschiedene Anfragen stellen, die aus unterschiedlichen Zielen und Problemen resultieren, sodass auch die Bildung eines Nutzermodells durch das IRS keinerlei Garantien besitzt, um Unsicherheiten bezüglich der Ziele des Nutzers zu verringern. Das gleiche Argument kann für Langzeit-Nutzer und periodische Nutzer etwa im Kontext von Information-Filter-Systemen angewendet werden, deren Anfragen, Ziele und Probleme einer zeitlichen Drift unterliegen kann.

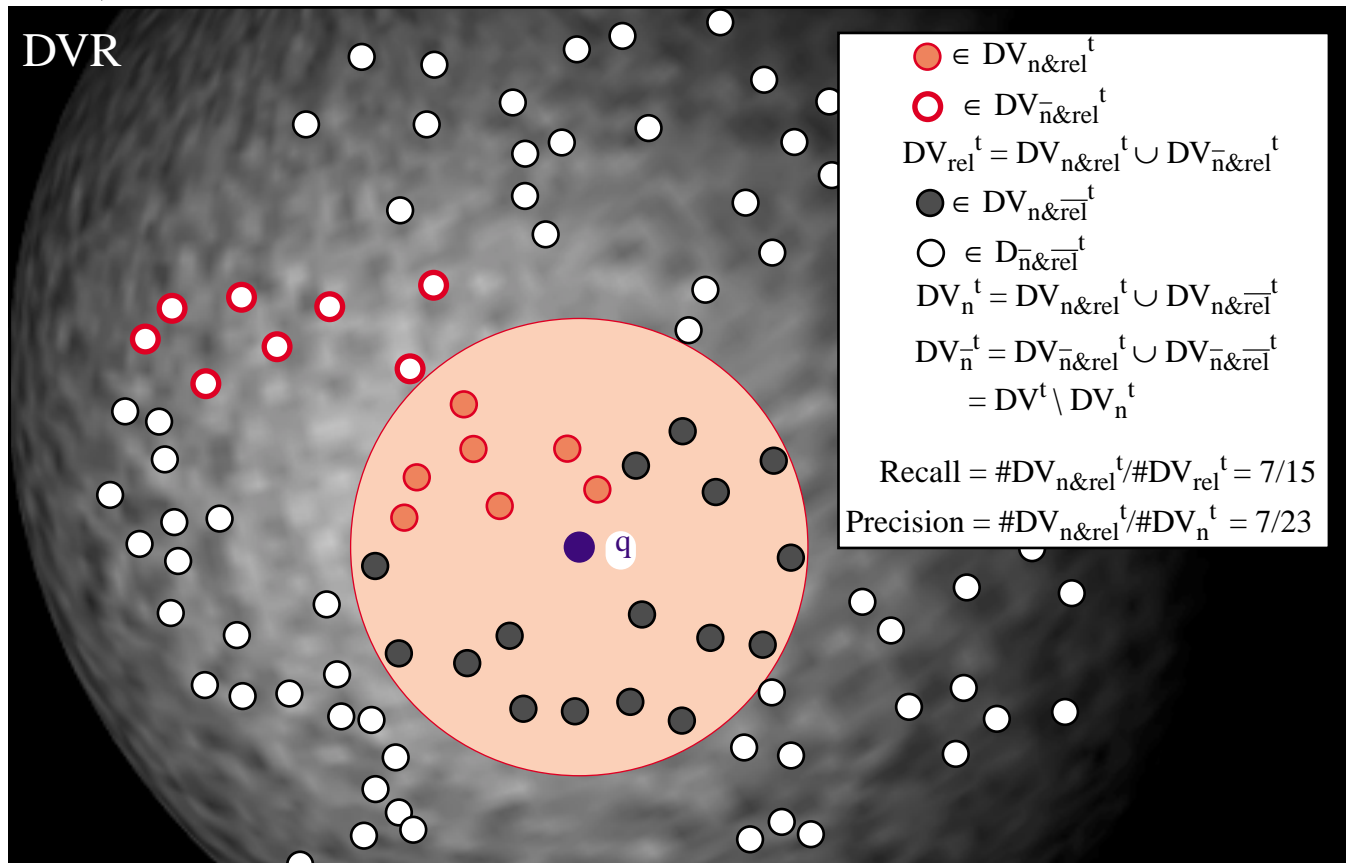
Die Bildung von Hypothesen über zukünftige Nutzer hat nicht nur Folgen für das Interfacedesign, sondern auch für die Wahl von Indexierungs- und Retrievalfunktionen. Der Versuch diese so zu wählen, dass einem „durchschnittlichen“ Agent geholfen wird, bestimmte Problemklassen zu lösen, ist eine sinnvolle Strategie, wenn von einer homogenen Agentenverteilung ausgegangen werden kann. Ist die Gesamtheit der Agenten jedoch bezüglich einer Vielzahl von Variablen divergent, und kann erwartet werden, dass das globale Leitbild der Individualisierung zu einer ständig größer werdenden Varianz in der Agentenverteilung führen wird, so ist die Auswahl einer einzelnen Indexierungs- und Retrievalfunktion und somit einer einzelnen Repräsentation von Informationsobjekten innerhalb des IRS keine geeignete Strategie. Diese Argumentation führt direkt zu dem Vorschlag der Polyrepräsentation in IRS (siehe Abschnitt 1.6) und Kapitel 3)).

### 1.2.6) Mehrziel-Anforderungen

Performancemessungen von IRS bezüglich einer Query werden durch Maße wie z.B. Recall als Quotient aus der Anzahl der nachgewiesenen relevanten Dokumente ( $D_{n\&rel}^t$ ) und der Anzahl der gesamten relevanten Dokumente ( $D_{rel}^t$ ) und durch Precision als Quotient aus der Anzahl der nachgewiesenen, relevanten Dokumente und der Anzahl der nachgewiesenen Dokumente ( $D_n^t$ ) bestimmt. In Abb. 1) werden die Dokumentvektoren anstatt der korrespondierenden Dokumente verwendet, um die benötigten Mengen und die Definition von Recall und Precision zu illustrieren. Recall und Precision stellen komplementäre Maße dar, d.h. die Maximierung des einen Maßes minimiert das andere Maß. Die daraus ableitbare Generalisierung, dass die Optimierung eines IRS ein Mehr-Ziel-Optimierungs-Problem (Rosenthal (1985[290]), Steuer (1986[324]), Ringuest (1992[281]), Dasgupta et al. (1999[90]), Veldhuizen (1999[347]), Zitzler (1999[375])) darstellt, wurde in der IR-Literatur jedoch nicht explizit dargestellt.

Die Mehr-Ziel-Sichtweise gilt insbesondere, wenn über Recall und Precision hinaus zusätzliche Performancemaße verwendet werden, die jeweils bestimmte Charakteristika von IRS beschreiben, und die teilweise Überlappungen untereinander bezogen auf die gemessenen Charakteristika besitzen. Erwähnt werden kann z.B. die Utility als Differenz der gewichteten Anzahl der relevanten Dokumente und der gewichteten Anzahl der nicht-relevanten Dokumente (siehe Vogt (1999: 11[352])). Weiterhin können Performancemaße verwendet werden, die auf Rang-Korrelationen basieren (Bartell (1994[28])) und Maße, die sich aus der Signal-Detection-Theory (Egan (1975[107]), Swets (1996[329])) ableiten (siehe Vogt (1999:11f[352])).

Abb. 1) Definition von Recall und Precision



Ohne eine Menge von Performancemaßen ergibt sich eine Mehr-Ziel-Anforderung an IRS direkt aus den restlichen fünf beschriebenen Problemstellungen. D.h. das Handling hochdimensionaler, nicht-linearer und dynamischer Zusammenhänge unter dem Vorliegen von Unsicherheit und Vagheit in einem IRS ist bereits eine Mehr-Ziel-Anforderung.

### 1.3) Methodentransfer

Die vorgestellten Einzelprobleme und ihre Kombination sind nicht ausschließlich für das IR spezifisch, d.h. es existieren andere Untersuchungsbereiche, in denen ähnliche Problemfelder auftauchen. Z.B. zeichnet sich ein Produktionsplanungssystem mit einer Menge unterschiedlicher Maschinen und einer zeitlich geordneten Liste an Jobs durch einen hochdimensionalen, diskreten Suchraum aus, in dem durch Schedulingverfahren nach optimalen bzw. suffizienten Maschinenbelegungsplänen gesucht werden muss

(Blazewicz et al. (1996[50])). Die damit verbundenen Probleme sind nicht-linear und nicht-stationär, wenn der Ausfall von Maschinen und der Ausfall bzw. die neue Aufnahme von Jobs disponiert werden müssen. Durch die Nicht-Beeinflussbarkeit dieser Faktoren durch das System ist es an die offene Welt gekoppelt, was Unsicherheiten und Vagheiten erzeugt. Durch zeitliche Constraints und monetäre Constraints entstehen zudem Mehr-Ziel-Anforderungen an die Leistung eines solchen Systems, sodass quasi alle aufgezeigten Problemfelder ebenfalls vorliegen.

Eine andere Domäne, in der Problemfelder dieser Art auftreten, ist die menschliche Kognition und ihre Modellierung durch kognitive Architekturen (Anderson (1990[8], 1991[9], 1993[10]), Newell (1990[233]), VanLehn (1991[345])). Dies dürfte nicht überraschen, da die Prozesse des IR, d.h. Indexierung und Retrieval, in Bibliotheken von Menschen lange vor der Verwendung von Computern durchgeführt wurden und immer noch durchgeführt werden. Wird das IR als Problem betrachtet, das von einem kognitiven System hinreichend gut gelöst werden kann, so liegt der Schluss nahe, kognitive Architekturen direkt als IRS zu verwenden. Dies ist jedoch mit einer Vielzahl von Problemen verbunden, die alle darin begründet sind, dass eine Menge von kognitiven Architekturen vorgeschlagen wurde, jedoch keine Architektur bezüglich Effektivität und Effizienz in den relevanten Problemstellungen dominierend ist. Doch selbst bei der Verwendung einer der Architekturen, die man bezogen auf Effektivitätskriterien als eine gute Architektur bezeichnen kann, ist ihre Leistungsfähigkeit gegenüber einer menschlichen Kognition immer noch bescheiden. Auf die kognitive Psychologie zu warten, um damit die Probleme des IR zu lösen, wäre eine vergleichbar verfehlte Strategie der IR-Forschung wie das Warten auf die KI in den 1980'er Jahren.

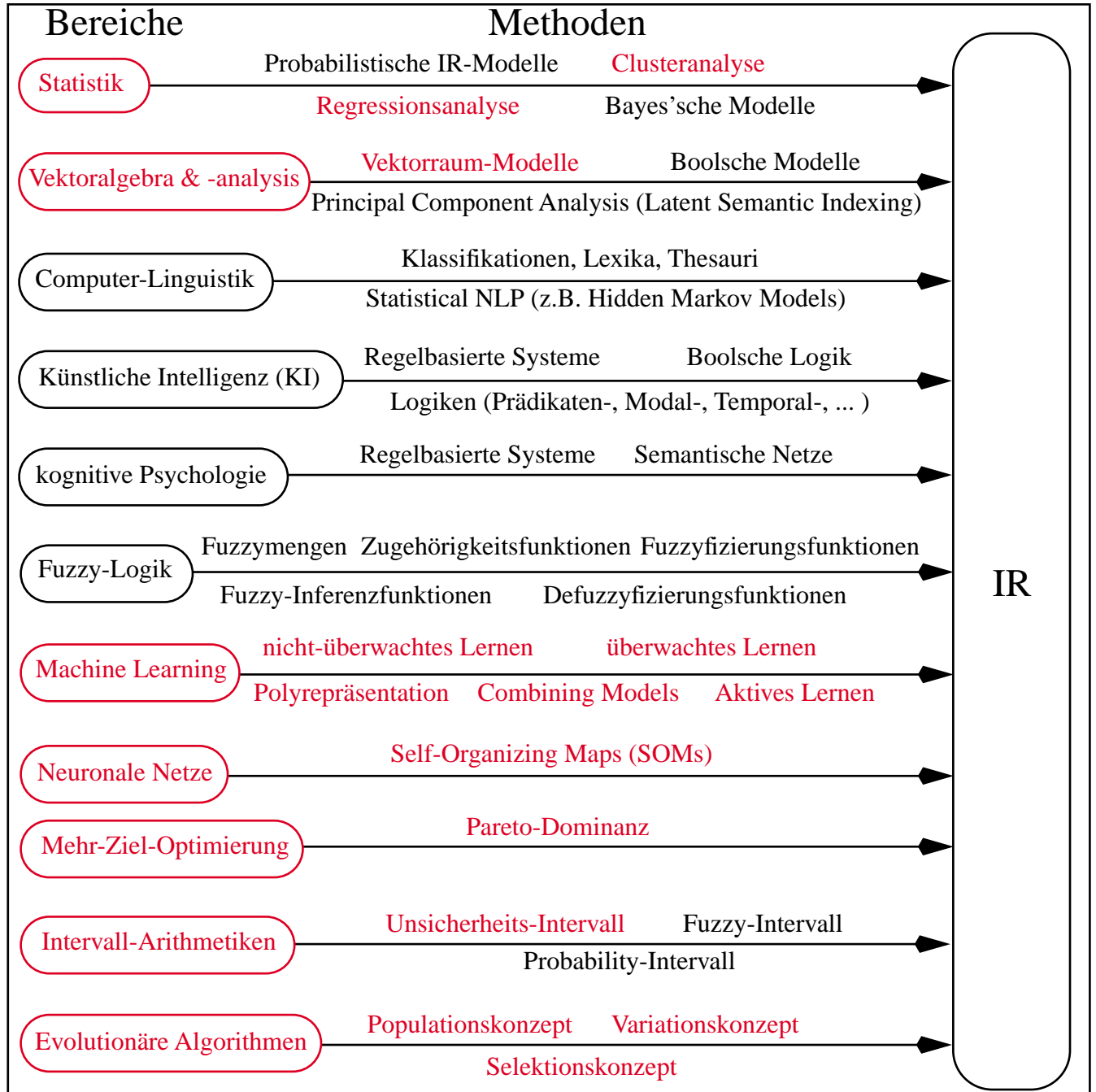
Es existiert eine Reihe weiterer Forschungsbereiche, welche mit gleichen bzw. ähnlichen Einzelproblemen konfrontiert sind, wie z.B. Machine Learning, Statistik, Optimierung und Control-Theory. In diesen Forschungsbereichen wurden eine Vielzahl von Methoden entwickelt, die zur Lösung der bereichsspezifischen Probleme eingesetzt werden. Durch die Erkenntnis, dass es sich um gleiche bzw. ähnliche Einzelprobleme handelt, wird ein Methodentransfer initiiert (siehe Abb. 2)), der über die gesamte Forschungshistorie des IR verfolgt werden kann. D.h. eine häufig beobachtbare Forschungslogik im IR besteht im Transfer von Lösungsmethoden aus Bereichen, in denen ähnliche Problemstellungen auftreten.

Betrachtet werden die Forschungsbereiche und die transferierten Methoden in Bezug auf die beschriebenen Problemfelder des IR, wobei Interdependenzen zwischen den Forschungsbereichen nicht detailliert beschrieben werden sollen wie z.B. Statistik und Computerlinguistik mit statistischen Natural-Language-Processing-Methoden (siehe jedoch Sprachlernen mit Evolutionären Algorithmen (Witten & Smith (1996[361]))).

Statistik ((Overbeck-Larisch & Dolejsky (1998[242]), Schuchmann & Sanns (1999a[303], b[304]), Abell & Braselton (1999[2])) beschreiben das Auftreten und die Eigenschaften von nicht-deterministischen Prozessen, d.h. es werden Unsicherheiten mit Hilfe von Konzepten wie z.B. Verteilungen bei der parametrischen Statistik beschrieben. Die direkte Methodenübertragung sind probabilistische Modelle des IR (Robertson (1977a[283], b[284], c[285], 1979[286]), Panyr (1986c[249]), Harter (1975[159]), Wong & Yao (1995[371])); Crestani et al. (1998[81])), bei denen z.B. für Merkmale die Wahrscheinlich-

keit berechnet werden, dass sie ein Dokument hinreichend gut beschreiben, oder bei denen für Dokumente die Wahrscheinlichkeit berechnet wird, dass sie bezogen auf eine Query relevant sind. Verbunden ist damit das „Probability Ranking Principle“ (Robertson (1977c[285])), bei dem die Dokumentergebnismenge in eine geordnete Liste überführt wird, indem als Ordnungsprinzip die fallende Wahrscheinlichkeit der Relevanz eines Dokumentes bezüglich der aktuellen Query verwendet wird.

Abb. 2) Methodentransfer ins Informations Retrieval



Die Zuordnung von Relevanzwahrscheinlichkeiten bei probabilistischen Modellen verwendet unterschiedliche Verfahren der schließenden Statistik wie Inferenz-Netze (Turtle & Croft (1990[341], 1991[342]), Turtle (1991[344]) oder Inferenzstrategien auf der Basis von Bayes-Verfahren (allg. siehe Martignon & Laskey (1995[211]); Turtle & Croft (1990[341], 1991[342]), Cooper (1991[78]), Kwok (1995[194]), Larkey & Croft (1996[196]), Henrion et al. (1997[161])).

Clusteranalysen (Eckes & Rossbach (1990[103]), Jain et al. (1999[175]), McLachlan & Basford (1987[214])) gehören in den Bereich der nicht-parametrischen Statistik (Schuchmann & Sanns (1999a[303])) und erzeugen aus einer Gesamtmenge von Objekten wie Dokumenten oder Merkmalen entsprechend der Repräsentation dieser Objekte Teilmengen (Cluster), die flach oder hierarchisch strukturiert werden können. Wie bereits beschrieben, werden Clustermethoden im IR im Kontext der hochdimensionalen Repräsentationen der Informationsobjekte eingesetzt (curse-of-dimensionality, Bellmann (1967), Warwick & Kárny (1997[357])).

Mit Regressionsmodellen (Ryan (1997[292]), Rawlings et al. (1998[273])) wird versucht, einen funktionalen Zusammenhang zwischen zwei oder mehr Variablen zu beschreiben, mit denen Objekte repräsentiert sind. Im Rahmen dieser Arbeit nimmt die Relevanz-Regression als Spezialfall der Relevanz-Approximation eine entscheidende Stellung ein (siehe Abschnitt 1.7) und Kapitel 4), siehe auch Cooper et al. (1992[79]); Gey (1994); Robertson & Walker (1994[287])). Die Zuordnung von Relevanzwerten zu Dokumenten bezüglich einer Query durch Regressionsmodelle ist eine alternative Vorgehensweise zu der Zuordnung von Relevanzwahrscheinlichkeiten in probabilistischen IR-Modellen, wobei Relevanzwerte bei Regressionsverfahren aus dem Intervall  $[0, 1]$  auch als Wahrscheinlichkeiten eines probabilistischen Modells interpretierbar sind.

Das Vektorraum-Modell des IR (Salton (1971[293]), Panyr (1987a[251])), auf das sich die vorliegende Arbeit bezieht, baut auf Begriffen wie Vektorraum, Distanzfunktion und Metrik auf, die aus der Vektoralgebra und Vektoranalysis (Strampp (1996[325]), Trostel (1993[339], 1997[340])) stammen. Verfahren des Latent-Semantic-Indexing (Dumais et al. (1988[100]), Furnas et al. (1988[136]), Bartell et al. (1992[27]), Soboroff et al. (1998[315]), Hofmann (1999[165]), Ando (2000[11])) basieren ebenso auf dem Vektorraummodell wie konnektionistische Verfahren des IR (Scholtes (1993[301]), Bachelier (1995[14]), Cunningham et al. (1997a[85], b[86])).

Boolsche Modelle des IR (Waller & Kraft (1979[354])) sind die ältesten Modellierungen, die Beziehungen zum Vektorraum-Modell sowie zu logischen IR-Modellen besitzen (Boolsche Logik). Grundlage von Queries sind Terme, die mit den Operatoren „Und“ ( $\wedge$ ), „Oder“ ( $\vee$ ), „Nicht“ ( $\neg$ ) verknüpft werden, wobei Klammern hinzu kommen können, um strukturierte, hierarchische Ausdrücke zu erzeugen. Dokumente werden im einfachsten Fall als Binärvektoren repräsentiert, d.h. jedem Term wird ein Komponent des Binärvektors zugeordnet, wobei die Komponente mit 1 belegt wird, wenn das Merkmal in dem zu repräsentierenden Informationsobjekt auftritt bzw. die Komponente wird mit 0 belegt, wenn das Merkmal nicht auftritt. Eine Retrievalfunktion prüft, ob ein Dokument bezüglich der Query passend ist, d.h. ob bestimmte Terme enthalten sind oder ob bestimmte Terme nicht auftreten. Nachgewiesene Dokumente können nicht unterschieden werden, sodass als Ergebnis eine Menge und keine nach einem Rankingkriterium geordnete Liste vorliegt. Erweiterte Boolsche Modelle versuchen ein Ranking zu erzeugen, indem sie Gewichte der Terme in den Dokumenten bzw. in der Query verwenden, wenn diese berechnet werden bzw. vom Agenten spezifiziert werden (Lee (1994[198]), Grossman & Frieder (1999:58ff[154])). Dies kennzeichnet einen Übergang zu vektorraumbasierten Modellen, da eine Distanzfunktion zwischen einer Query und einem Dokument auf der Basis der Terme in der Query und den Gewichten als Distanz in einem metrischen Vektorraum betrachtet werden kann.

Computer-Linguistik allgemein versucht die Spracherkennung, das Sprachverstehen (Ram & Moorman (1999[270])) und die Spracherzeugung beim Menschen zu modellieren und zu implementieren (siehe z.B. Batori et al. (1989[29])), wobei eine Vielzahl symbolischer sowie subsymbolischer, d.h. konnektionistischer Ansätze (z.B. Scholtes (1993[301])) vorgeschlagen wurden. Zu den symbolischen Ansätzen gehören die Verwendung von Klassifikationen, Lexika und Thesauri, während statistische Methoden der Verarbeitung natürlicher Sprache (Statistical Natural Language Processing (Manning & Schütze (1999[210])) die Schnittmenge zwischen den symbolischen und konnektionistischen Ansätzen darstellen. Zu den statistischen Verfahren gehören insbesondere Hidden Markov Modelle (Buchholz (1991[60]), Sinclair (1992[311]), Davis (1993[92]), Buchholz et al. (1994[61]), Elliott et al. (1995[109])), für die auch Anwendungen in IR existieren (Miller et al. (1999[219])). Da in dieser Arbeit symbolische Ansätze wie die Verwendung von Lexika und Thesauri nicht betrachtet werden, besteht ein Anknüpfungspunkt bei den statistischen und konnektionistischen Verfahren, wobei insbesondere der Ansatz zu erwähnen ist, dass Dokumente, Queries und Merkmale (Features) als Zeichenketten eines endlichen Alphabetes betrachtet werden. Etwas periphratisch werden im Rahmen der Polyrepräsentation Sprachgenerierungsmodelle betrachtet, die aus einem Dokument oder einer Query eine Menge von Dokumenten bzw. Queries erzeugen, die nachfolgend indexiert werden. Dabei werden insbesondere Zeichenketten-Resampling-Verfahren betrachtet wie das Moving-Blocks-Verfahren (siehe Abschnitt 2.2.3)).

Künstliche Intelligenz (KI) (Görz (1993[151]), Bundy (1997[64])) versucht allgemein Hard- und Software zu entwickeln, die Leistungen erbringen, für welche beim Menschen das Attribut intelligent verwendet würde. Dabei wird keine biologische oder psychologische Plausibilität wie bei Neuronalen Netzen oder Modellen der kognitiven Psychologie gefordert, sondern es werden ausschließlich Effektivität und Effizienz als Performancemaße verwendet. Die KI hat seit Mitte der 50'er Jahre eine große Anzahl von Repräsentationsformen und Lernformen entwickelt, von denen einige im Kontext des IR eingesetzt wurden, wobei regelbasierte Systeme und logikbasierte Systeme angesprochen werden sollen.

Regelbasierte Systeme (Produktionssysteme; Newell (1990[233])) besitzen eine Wissensbasis in Form einer Menge von Wenn-Dann-Regeln. Kombiniert wird diese Wissensbasis mit einem Regelinterpreter, der in Abhängigkeit von einer Menge momentan aktiver Wissens Elemente in einem Arbeitsspeicher (Working-Memory) ermittelt, welche Regeln anwendbar sind, indem er die Wissens Elemente mit den Wenn-Teilen (If-Teile) der Regeln vergleicht (matched). Je nach der Art des Interpreters führt er alle anwendbaren Regeln aus, eine Teilmenge oder die erste passende Regel, die er findet. Der Dann-Teil (Then-Teile) der Regeln kann Wissens Elemente im Arbeitsspeicher modifizieren, löschen, neue erzeugen, externe Aktionen auslösen u.a. Die Regeln in der Wissensbasis werden meist von Menschen erzeugt, bzw. es werden Lernmechanismen implementiert, die neue Regeln erzeugen können, was den Bereich des Machine Learnings berührt.

Der Einsatz von regelbasierten Systemen im IR überschneidet sich mit der Argumentation des Einsatzes von Modellen der kognitiven Psychologie im IR. Die Hypothese der kognitiven Psychologie, dass kognitive Prozesse durch Produktionssysteme adäquat modellierbar sind (Anderson (1990[8], 1991[9], 1993[10]), Newell (1990[233])), ist der Grund, warum diese Systeme im IR eingesetzt werden. Die Indexierung wurde und wird von menschlichen Experten durchgeführt, bevor Computersysteme sich an dieser Aufgabe versuchten. Ziel ist es in diesem Kontext, die Indexierungsregeln zu modellieren, von denen