

Zusammenfassung IRG

René Bernhardsgrütter, 13.05.2014/16.06.2014
IB=Informationsbedürfnis.

Grundlagen

Information Retrieval

Viele un-/strukturierte Informationen organisieren und auffindbar machen.

Pyramide

#3: Daten: Fakten in codierten Form. Tragen per-se keine Bedeutung, können aber korrekt od. falsch sein.

#2: Information: Daten+Bedeutung. Information ist i-/relevant. Benötigt, um Aufgaben zu erledigen, für Aufgabe mehr/weniger relevant.

#1: Wissen: Verarbeitete, vernetzte Information (Erkenntnisse). Z.B. Bestellung abwickeln. Häufig interne und externe Informationen vernetzen.

Retrievalproblem

„Auffinden von möglichst viel relevanter Information bei gleichzeitigem Minimieren von ebenfalls gelieferten irrelevanten Information.“

Nicht nur Informationen wieder finden, sondern vor allem neue Information finden. Information wird indirekt geliefert, in Form von „relevanten“ Dokumenten.

Folge: Ein perfektes Retrievalresultat losgelöst von Benutzer und Kontext gibt es nicht.

Sprachproblem

- Sprache ist nicht „eindeutig“: Synonyme (eine Bedeutung – mehrere Wörter), Homonyme (mehrere Bedeutungen – ein Wort) Umschreibungen, Metaphern, Wortformen (Singular, Plural, Fälle, etc.).
- IB ungenügend verbalisiert und formuliert.
- Dok./Infos im System unstrukturiert/ inhomogen.
- Irreführender Inhalt.
- Autorität, Quelle, Aktualität, Urheberrecht, Einsammeln der Dokumente.

Relevanz

Immer subjektiv: Vor-/Hintergrundwissen, Reihenfolge des Auffindens, wandelnde IB, persönliche Präferenzen, Vollständigkeit der Antwort.

Konsequenzen

Unschärfer Relevanzbegriff führt zu wskbasierter Lösung:

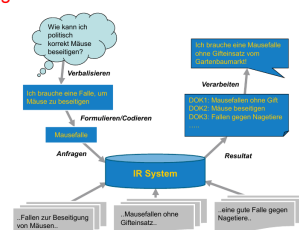
- Für User vermutlich relevante Dokumente top gelistet.
- Resultat fast nie vollständig „korrekt“: relevante Dokumente fehlen, od. zusätzlich irrelevante gefunden.
- **Scharfe Kriterien (ja/nein) sind ungeeignet**, da User Anz. und Form der gesuchten Dokumente der Anfrage kennen müsste, welche das gewünschte Resultat liefert -> Paradox.
- Gute Retrieval-Systeme erlauben alles zu formulieren was bekannt, ohne zu viel od. zu wenig zu finden.

IR-Prozess

Verbalisierung: „Verstehen des Problems“: richtige Begriffe, Vollständigkeit. => **Problem muss verstanden werden, um es richtig zu verbalisieren.**

Codierung:

„Verstehen des Systems“: richtige Operatoren, etc..
=> Die Resultate müssen bekannt sein, um Anfrage richtig zu codieren.



Das Resultat ist immer nur so gut wie die Anfrage!

Suchparadoxon

Google kann sich "Vereinfachungen" erlauben dank dem Suchparadox: **Es ist einfacher, in mehreren Milliarden Dokumenten zu suchen als in mehreren Tausend.**

- Grosse Datenmengen => mehr Redundanz.
- Information mit "beliebigen" Verbalisierungen findbar.
- Benutzer von Google sind häufig präzisionsorientiert, wenige gute Treffer reichen.

Auf "kleinen" Datenmengen ist Übereinstimmung zwischen Informationsbedürfnis und Dokumenten schwerer nachzuweisen.

Iterative Suche

Verständnis des Benutzers für Informationsbedürfnis ändert sich mit gesammelter Information => iterativ:

- Unsterstützung bei Umformulierung: automatische Erweiterung der Anfrage, Suche nach ähnlichen Dokumenten.
- Werden relevante Dokumente gefunden, ändert sich das Verständnis des IBs.
- Benutzer kann beim Verständnis unterstützt werden, z. B. durch automatische Kontextanalyse.

Information Retrieval Paradigmen

Pull

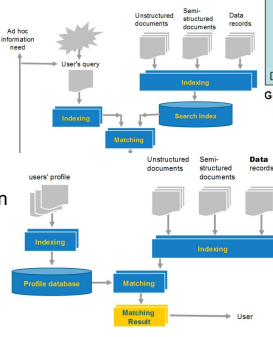
Aufgrund eines Ad-hoc-IBs Dokumente mit relevanten Informationen suchen.

Bsp: Ad-hoc-Suchen wie von Bundesgericht.

Push

Dokumente mit relevanten Informationen aus einem Dokumentenstrom herausfiltern und weiterleiten.

Bsp: RSS Reader (Bringdienste)



Browse

Neue Dokumente kategorisieren und einordnen.

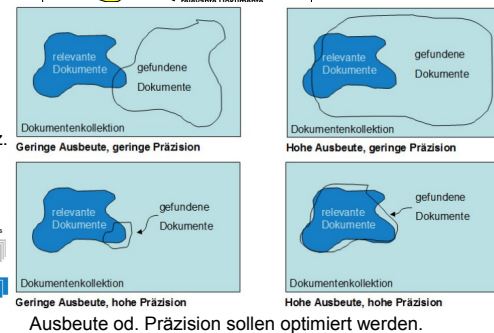
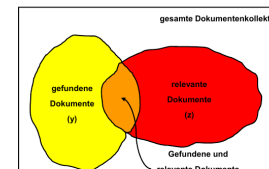
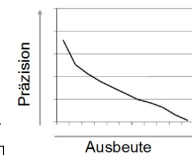
Bsp: dir.yahoo.com

IR vs. DB-Suche

- Datenbanken liefern für strukturierte Information mit kontrolliertem Vokabular **perfekte** Resultate.
- Daten in Datenbanken sollten unabhängig sein von Applikation. Redundanz wird vermieden.
- Elemente sind entweder Teil der Resultatmenge od. nicht (binäre Unterscheidung).
- Boole'sche Kriterien zur Selektion Geeignet für die Suche in (hoch-)strukturierter Information mit kontrolliertem Vokabular.

Qualität von Retrievaleffekt

Ausbeute und Präzision modellieren Annahme, dass möglichst viel relevante, und möglichst wenig irrelevante Information gefunden werden soll.



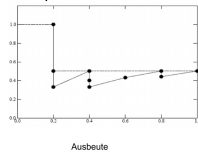
$$Präzision := \frac{\# \text{relevante Dokumente im Resultat}}{\# \text{Dokumente im Resultat}}$$
$$Ausbeute := \frac{\# \text{relevante Dokumente im Resultat}}{\# \text{relevante Dokumente in der Kollektion}}$$

Da beide Masse mengenbasiert sind, widersprechen sie sich oft: hohe Ausbeute -> geringe Präzision, hohe Präzision -> geringe Ausbeute

Beispielauswertung

■ Berechnung von Ausbeute und Präzision auf Ranglisten

Rang	Relevant?	Ausbeute	Präzision	Prz. interpoliert
1	+	0.20	1.00	1.00
2	-	0.20	0.50	0.50
3	-	0.20	0.33	0.50
4	+	0.40	0.50	0.50
5	-	0.40	0.40	0.50
6	-	0.40	0.33	0.50
7	+	0.60	0.43	0.50
8	+	0.80	0.50	0.50
9	-	0.80	0.44	0.50
10	+	1.00	0.50	0.50



Probability Ranking Principle (PRP)

Resultat sortiert nach Relevanz-WSK. **Optimal** unter Berücksichtigung sämtlicher zur Verfügung stehender Informationen und geeigneter Annahmen. Ist eher Hypothese als Prinzip.

PRP halt folgende math. Eigenschaften:

- Präzision an beliebigem "cut-off point" wird optimiert.
- Ausbeute an beliebigem "cut-off point" wird optimiert.
- Kosten der Auswertung (relevant=positiv, irrelevant=negativ) werden optimiert.

=> WSK-basierte Ranglisten theoretisch fundiert

Problem

q = "Terrorismus, bekämpfen" = $\phi_1 \phi_2$

D1 = "Gegenmassnahmen gegen Terrorismus" = $\phi_3 \phi_1$

D2 = "Kampf gegen den Terror" = $\phi_2 \phi_1$

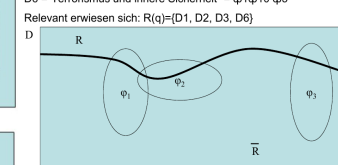
D3 = "Sicherheit bei asymmetrischer Bedrohung und asymmetrische Sicherheit" = $\phi_5 \phi_6 \phi_7 \phi_6 \phi_5$

D4 = "Terror bekämpfen" = $\phi_1 \phi_2$

D5 = "Extremismus und Gewalt" = $\phi_8 \phi_9$

D6 = "Terrorismus und innere Sicherheit" = $\phi_1 \phi_{10} \phi_5$

Relevant erwiesen sich: $R(q) = \{D1, D2, D3, D6\}$



Indexierung und Vergleichen

Term/Merkmal: eindeutiges Wort

Token: Auftreten eines Terms

Wort: Einheit zwischen Trennzeichen bei Tokenisierung.

Merkmals Häufigkeit: Anz. Vorkommen von Merkmalen in Dokument.

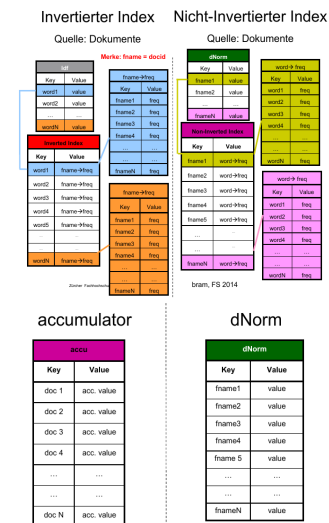
Dokumenten Häufigkeit: Anz. Dokumente, die Begriff enthalten.

Inhaltstragende Wörter

Manche Wörter tragen wenig *lexikalischen* Inhalt. Z.B. Artikel, Partikel, Pronomen, Konjunktionen, ..

Aber: Bedeutung von "nicht" ist wichtig! Für Retrieval interessant sind Wörter, welche inhaltstragend sind und Text "auszeichnen".

MiniRetrieve



Der Akkumulator summiert das Produkt von idf und Termhäufigkeit für jedes Wort, das im entsprechenden Dokument vorkommt.

Dokumentenorm „dNorm“ wird für alle Dokumente vorberechnet

Der idf (Inverse document frequency) wird für alle Wörter in allen Dokumenten vorberechnet, ggf. auch für Anfragebegriffe mit df=0.

Evaluation

Nötig, um „Leistung“ des Systems zu bewerten.

Zutaten einer Evaluation

Aufgabenstellung: Motivation, Ziele? Zielgruppe? Interne od. externe Evaluation? Black/Whitebox? Was ist Benchmarkgrösse?

Ausgestaltung: Was für Leistungsfaktoren (Umgebungsvariablen, Systemparameter)? Was für Leistungsmaassstäbe? Was für Leistungsmasse (Effektivität, Effizienz, Akzeptanz)?

Testdaten: Daten über Users, IB, Dokumente mit Relevanz bzgl. IB.

IR-Evaluation nach Cranfield/SMART

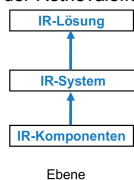
Heute meistverbreitet in akademischen Kreisen.

3 Stufen der Evaluation der Retrievaleffektivität:

■ Ganzer Wissensbeschaffungsprozess (UI-Fragen, etc)

■ Ganzes Retrievalsystem (Anfrage -> Dokument)

■ Komponenten des Retrievalsystems (Stemmer, etc.)



Ebene

Aufgabenstellung

• Retrievaleffektivität des IR-Systems soll eval. werden.

• „Interne“, direkte Evaluation angestrebt, d. h., Effektivität direkt messen, nicht über Anteil an

grösserem Resultat.

- Evaluation soll experimentell erfolgen.
- System während Evaluation als „black box“ behandelt.
- Evaluation erfolgt im Vergleich zu optimalem Resultat.
- Soll quantitative Evaluation durchgeführt (Effektivität).

Ausgestaltung: Labortest

- Setup von operationellen Umgebungen abgekoppelt.
- wird von spezifischen Benutzern und Interpretationen ihrer IB abstrahiert (-> Anz. Umgebungsvariablen minimiert).
- Leistung als durchschnittliche Leistung über Anz. von Retrievalvorgängen gemessen.
- Leistungsmasse sind Präzision und Ausbeute.
- Dokumentdaten geeignet bestimmt und „eingefroren“.
- „Batch-Setup“ und Experimente beliebig **wiederholbar**.

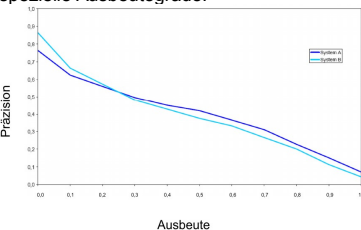
Testdaten: Testkollektion

- Reihe IBn von fiktiven Benutzern.
- „eingefrorene“ Menge Dokumenten als Suchdaten.
- Relevanzbeurteilungen von Dokumenten relativ zu IB.

Average Precision

non-interpolated: Durchschnitt der Präzisionswerte an Rängen aller relevanten Dokumente.

interpolated: Durchschnitt der Präzisionswerte für spezielle Ausbeutegrade.



Beliebteste „Ein-Zahl-Mass“. Konvergiert gegen Fläche unter dem *Ausbeute/Präzisions-Graph*.

Aber: Widersprüchlichkeit von Ausbeute und Präzision bedingt, dass „Average Precision“ nicht immer ein praxisrelevantes Mass darstellt.

Mean Average Precision

Durchschnittliche Average Precision über alle Anfragen.

Ausbeute

Um Ausbeute zu bestimmen, sämtliche Dokumente nach Relevanz bewerten. Heute aber unrealistisch.

=> Vorhandene Testkollektionen wenn möglich nutzen.

=> Ausbeute auf alternative Art bestimmen.

Evaluations-Kampagnen

Idee: wenn genügend viele verschiedene Teilnehmer mit unterschiedlichen Systemen Evaluation durchführen, werden (fast) alle relevanten Dokumente von mindestens 1 System gefunden – müssen nicht alle Dokumente gelesen werden, um Ausbeute zu bestimmen.

Konkret z. B. als TREC-Kampagnen.

Gossip

Bei 1 Feedback eines unzufriedenen Kunden haben 20 andere bereits selbes Problem gehabt, geben aber kein Feedback (20 mal Faust-im-Sack)

- System darf sich keine groben Ausrutscher erlauben.
- Avg. Performance kompensiert nicht für Ausrutscher.
- Man sollte sich nicht allzu viele Gedanken über **positive** Meldungen machen.

Known Item Retrieval

Idee: mit Suchmaschine nach „bekannten“ Dokumenten suchen, simuliert „Da war doch was“.

Mean Reciprocal Rank (MRR) = 1/R: Effektivität.

Ermittelt Durchschnitt über Anz. von Anfragen.

Kategorisierung und Klassierung

Kategorie: Dokument kann in keiner, einer od. mehreren Kategorien sein.

Klasse: Jedes Dokument in genau einer Klasse (es gibt möglicherweise auch Klasse *unklassiert* od. *andere*).

Dokumentkategorisierung

Kategorien bieten Mehrwert bei (manueller) Indexierung, Suche nach Information in Dokumentenarchiven od. News Feeds.

Indexierung: Verschlagwortung

Routing/Filtering: Themenprofil erstellen

Clustering: Gruppieren von Kollekt. (Memos, E-Mails,...)

Annotation: Gruppieren von Dokumenten

Informationsstrukturen

Bei *Informationsstrukturen* geht es immer darum, Dokumente/Informationsobjekte in Verbindung setzen.

Typische Informationsstrukturen: Thesauri, Klassifikationen, Prozess-Modelle (Geschäftsprozesse), UML, Inferenzmodelle aus der künstlichen Intelligenz, Navigationsstrukturen (Intranet), ...

Rocchio Modell

Modelliert Kategorie C mittels Repräsentanten c. Dieser ist wie Dokument ein Vektor, aber hypothetisch. Initial aus positiven Beispielen generiert, regelmässig aktualisiert.

Algorithmus: Neue Dokumente d als Anfrage zum Vergleich mit den Repräsentanten. Falls $s(d, c) > \Delta$ Dokument d der Kategorie c zugeteilt. Δ ist Grenzwert, geeignet zu bestimmen.

Vorteile: einfach implementierbar, extrem effizient.

Nachteile: Nicht robust wenn Anz. von negativen Instanzen gross. Festlegung von Parametern knifflig.

k-Nearest Neighbor (kNN)

Verwendet Ähnlichkeitsmass (Eukl. Dist., cos) und Regel, wie Dokumente in D Kategorien zuzuordnen.

Algorithmus: Bestimme k ähnliche Dokumente zu D, d.h. k nächsten "Nachbarn". Ordne D min. 1 Kategorien

von Nachbar zu.

Nachteil: Testkollektion nötig um Vorgang zu starten.

Erweiterte kNN-Klassifikation

Je weiter Dokument D vom Nachbar D_i entfernt (ϕ), desto weniger trägt es zum Entscheid bei, Dokument D in Kategorie C_i zuzuordnen.

Wobei $kNN(D)$ Menge von k nächsten Nachbarn von D. $a_{ij}=1$ falls Dokument D_i zu C_i gehört, $a_{ij}=0$ sonst.

Probleme: Richtige Wahl von k, Funktion sc, und max.

Anz. zugeordneter Kategorien. Auch Schwellwert ab wann ein Dok. in mehrere Kategorien.

Vorteile: Effektiv und einfach, stabile Schwellw. zu finden

Nachteile: Langsam, Wahl einzelnen Wertes "k" schwer.

Bayes Klassifizierung

Gegeben: Kategorie C_i mit angemessenen Anz. von bereits zugeordneten Objekten (Trainingsdaten).

Methode: Bilde statistische Modelle aus Kategorien, um bestimmen zu welcher Klasse neues Objekt D gehört. $P(t|C_i) \forall t$ bekannt, C_i , aber interessiert an $P(C_i|t)$ od. $P(C_i|D)$. D ist Menge von Merkmalen in Objekt/Dok. D.

Wsk, dass D zu C_i gehört:

$$P(C_i | D) = \frac{P(D | C_i)P(C_i)}{P(D)}$$

Mit $D=(t_1, M., t_n)$ und in Beziehung zu Klasse C_i :

$$P(D | C_i) = \prod_{j=1}^n P(t_j | C_i)$$

Gibt verschiedenen Wege um $P(t|C_i)$ zu berechnen: zähle Anz. Merkmale, binär (Nicht-/Vorkommen), gewichtet...

Vorteile: Sauberes Modell, einfach zu implementieren.

Nachteile: Performt sehr schlecht, typischerweise schlechter als andere einfache Verfahren.

Unabhängigkeitsannahme wohl zu simpel.

Regelbasierte Methode

Gewünschte Kategorie mittels Regeln beschreiben.

Bsp: Selektion geeigneter Beispieldokumenten, was für Suchanfrage ergibt diese Dokumente als Resultat?

Problem: Extrem schwierig, beständige (lange gültige)

Anfragen zu Aufwand Konsistenz (manuell konstruierte Thesauri sind selten zyklentfrei) Unschärfe?

(Transitivität?) Vollständigkeit? Bedeutungen abhängig von Zeit, Betrachter, thematischen Kontext. Wie global kann Thesaurus sein? Thesaurus kann helfen, Begriffe konsistent zu verwenden.

Use-Case: Regelbasierte Kategorisierung gut für "scharfe" Konzepte => Expertensysteme.

Bsp: Regel für Kategorisierung von AUSTRALIAN DOLLAR sieht z.B. so aus: [australian-dollar-concept] or ([dollar-concept] and [australia-concept] and not [us-dollar-concept] and not [singapore-dollar-concept])

Verbesserungen dann möglich, wenn zu suchenden Konzepte in Feldern vom Text auftreten (z.B. Titel).

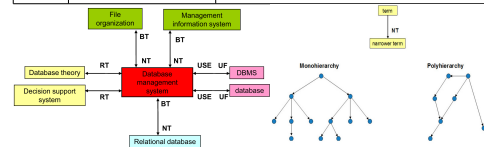
Vordefinierte Informationsstrukturen

Traditionelle Verfahren um Beziehungen zwischen Konzepten/Kategorien, Worten und Phrasen zu definieren. Unterscheidet zw. **Deskriptoren** und **Nicht-Deskriptoren**, so genannte „lead-in terms“. Klassifikationen bestehen nur aus NT-Relationen. Sind eher statisch und hauptsächlich im Bibliotheksumfeld zu finden.

Use-Case: Indexierungszwecke, Anfrageformulierung, verhindert „mismatch“ zw. Anfrage- und Dokumentmerkmalen. Abstimmung eines Retrieval-Alg.

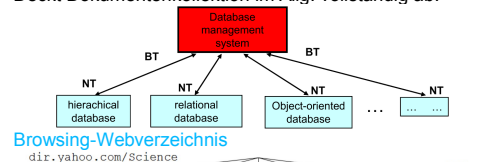
Komponenten: kontrolliertem Vokabular (z.B. Kategorien!), Beziehung zw. einzelnen Merkmalen des Vokabulars, und Info zu Merkmalen und Beziehungen.

USE	synonym	use another term as descriptor
UF	used for	use this term instead
NT	narrower term	restricted concept
BT	broader term	umbrella concept
TT	top term	concept at head of hierarchy
RT	related term	similar concept

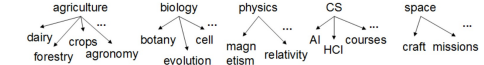


Herausforderungen: Aufwand, Konsistenz, Unschärfe, Vollständigkeit, Bedeutungen abhängig von Zeit, Betrachter, thematischen Kontext.

Klassifikation: Wenn Thesaurus nur NT (+BT) Beziehungen hat. Handelt dann um Monohierarchie. Deckt Dokumentenkollektion im Allg. vollständig ab.



Browsing-Webverzeichnis



Herausforderungen: Verwaltungsaufwand, Vollständigkeit, Bedeutungen ändern, Klassifikation muss ändern können, Menschen klassifizieren anders, etc..

Clustering

Menge von Objekten in "Cluster" von ähnlichen Objekten. Cluster können beliebige Untercluster besitzen.

Bsp: Resultat einer Suchanfrage durch Clustering geordnet. Hift User, effizienter weiterzusuchen od. Überblick über Kontext zu erhalten. Navigationshilfe.

Social Tagging

Social Tagging beschreibt Prozess bei welchem Benutzer Metadaten in Form von Schlüsselwörtern hinzufügen und diese Daten mit anderen Nutzern teilen.

Tag: "Aufkleber", beschreibt Informationsobjekt näher.

Tagging: Tags zu Informationsobjekt adden u. ggf. teilen.

Eigenschaften: Spontan durch Community, keine fixen Regeln, kein kontrolliertes Vokabular, keine Hierarchie, keine Beziehung zwischen Tags, Vergleiche streben nach Konsistenz bei Kategorisierung, analog zum Verhältnis Verschlagwortung vs. Volltextindexierung.

Ressourcen: Dokumente, Artikel, Einträge, etc.. Linking zw. Ressourcen gut untersucht (PageRank/HITS).

Vorteile für Betreiber: Nutzer kategorisieren, Tags reflektieren Benutzervokabular. Users erstellen Metadaten. Flexibel auf Veränderungen, bessere Abdeckung als mit kontrolliertem Vokabular, kostengünstig, stärkere Bindung an Angebot.

Vorteile für Benutzer: Auffinden eigener Informationsobj., leichteres Auffinden von gem. Informationsobjekten, mehrere Benutzer teilen, Anreiz schaffen für andere, Selbstdarstellung, Ausdruck eigener Meinung.

Nachteile: Zersplitterung der Kat. da unkontrolliertes Vokabular, Singular/Plural, mangelnde Struktur der Metadaten, mangelnde Präzision, ungenaue Deskriptoren

Wortwolke (Tag Cloud)

Methode zur Informationsvisualisierung, Wortliste alphabetisch aufgeführt und Grösse eines Worts seiner Wichtigkeit entspricht.

Web Search

Daten im Web sehr unbeständig (40% monatliche Änderung), unstrukturiert, teilweise schlechte Qualität.

Ziel: Finde qualitativ hoch stehende Resultate (nicht unbedingt Dokumente!), die für den Benutzer relevant.

Resultat: Statische Seiten (Dokumente, z.B. Texte, mp3, Photos, Videos), Resultat nicht ausbeuteorientiert, ggf. Probleme mit dynamischen Seiten.

Bedürfnisse

informationell: etwas lernen (~40%) z.B. Was ist ISIS?

navigational: will zu Webseite (~25%) z.B. SBB

transaktional: man will etwas tun (~35%), z. B. Zugriff auf Service "Wetter in Zürich", Downloads, Shop, iTunes.

Graubereiche: Finde guten Hub, z.B. "Automiete in Seattle", Erkundungssuche "see what's there".

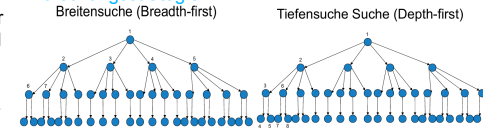
Spider

in-degree: Anz. Links, die auf Page zeigen.

out-degree: Anz. Links, die von Page weg zeigen.

Starte mit umfassenden Menge von URLs, von welchen Suche gestartet wird. Speichere Dokumente in D und Hyperlinks in E. Während Crawling wird Liste Q von URLs intern unterhalten. Wir extrahieren eventuell URL von Q od. fügen URL Q hinzu (Funktionen Dequeue() und Enqueue()).

Erforschungsstrategien



Tiefensuche erfordert Speicher nur für die Tiefe (d) mal den Verzweigungsgrad (b) ($O(bd)$) => Algorithmus benötigt aber zu viel Zeit, um nur 1 Zweig nachzugehen. **Breitensuche** erforscht von Rootwebseite gleichmässig nach aussen. Erfordert Speicher für alle Knoten des Graphen von früheren Ebenen ($O(b^d)$). Kompromiss nötig. Wie neue Links zu Q hinzugefügt (Enqueue()) und extrahiert (Dequeue()) hängt von Suchstrategie ab.

FIFO, first in first out: Ende von Q angeh., Breitensuche.

LIFO, last in first out: Anfang von Q angeh., Tiefensuche. Heuristisches Anordnen von Q ergibt fokussierten Crawler, der Suche auf interessante Seiten ausrichtet. Z. B. auf Sites, Ordner, Sprachen, etc. begrenzzbar.

Spreading Activation (SA)

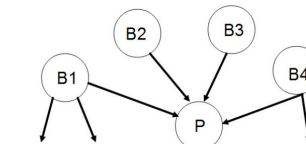
Grad der Ähnlichkeit zw. Di und Q, (bzw mit $\text{Sim}(Di, Q)$), wird durch bestimmte Anz. von Kreisläufen (normalerweise 1) zu verlinkten Dokumenten propagiert. Dazu Faktor λ verwenden. Nur Top r Dok. angepasst.

PageRank

Zu Beginn ist User auf zufälliger Page. Jeder Schritt bewegt User mit Wsk d zu zufällig gewählten Page (z. B. $d = 0.15$) od. zu zufällig gewähltem Nachfolger der aktuellen Page mit der Wsk $1-d$ (z.B. $1 - 0.15 = 0.85$).

PageRank einer Webseite: Wsk, dass User zu beliebigem Zeitpunkt auf Page ist.

Qualität von P: $Q(P) = Q(B1)/3 + Q(B2) + Q(B3) + Q(B4)/2$



$$PR^{c+1}(D_i) = (1-d)\frac{1}{n} + d \left[\frac{PR^c(D_1)}{C(D_1)} + \dots + \frac{PR^c(D_m)}{C(D_m)} \right]$$

$C(D_1)$: Anzahl Outlinks von D1. Faktor $1/n$ ist umstritten.

Eigenschaften: Berechnung rekursiv, offline, anfangs ist Page-Rank aller Knoten 1.0, konvergiert gegen "Random Surfer"-WSK, anfällig auf Spam.

Nachteil: Muss theoretisch nach jeder Änderung neu berechnet werden.

Kleinberg Model/HITS

HITS: Hypertext Induced Topic Selection. Site ist Hub (verweist auf Quellen) od. Autorität (besitzt Inhalt). Werte iterativ berechnen (nach 5 Its schon relativ stabil). Nach jeder It normalisiert => Summe der Quadrate = 1.

Link-Analyse/Bowtie Model

Betrachtet Anz. Pages einer Site und Anz. Links einer Site oder Page.

$$P(n) = b \cdot n^{-\beta}$$

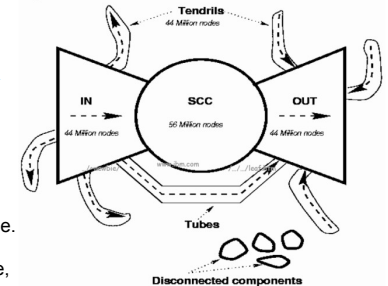
with $b > 0$ and $\beta > 1$

Anz. von Pages einer Site wird als n angenommen.

$P(n)$ ist die Wsk, dass die Site n Pages hat.

Enterprise IR

Suchsysteme für (un-)strukturierte Daten in Unternehmen; unterstützt wissensintensive Geschäftsprozesse.



Use-Case: Interne, vertikale

Kollektionen, firmenübergreifende Suchsysteme (integrative Funktion!), Intranets, Desktop Search, Extranets, spezielle Suchapplikationen.

Impact: Grossen Markt für Enterprise Search. 70% aller Anfragen direkt bei Sites.

Aufsetzen

- Problem verstanden? (Informationen in 50 Dokumenten suchen != Informationen 50 Mio. Dokumenten suchen).
- Wer sind die Benutzer, was für IB?
- Wer "besitzt" die Daten? Welche Dokumente werden erschlossen?
- Risiko: Es besteht die Gefahr, dass Bedürfnisse mit einem bestimmten Produkt od. einen bestimmten Technologie assoziiert werden.

Anforderungen

- geschäftrelevanten Entitäten (Produkte, Kundenbeziehungen, Lieferanten, etc.) „kennen“
- wichtigsten Benutzer- und Bedürfniskategorien „kennen“
- relevante Informationsquellen für welchen Prozess
- Transformation von unstrukturierter Information in strukturiertes Wissen unterstützen.

Technologien

- Topic Detection & Tracking
- Kategorisierung
- Clustering
- Information Extraktion: Message Understanding, Named Entity Recognition (erkennen von Personennamen, Ortsnamen, Produktnamen etc.)

Merkmal	Web IR	Enterprise IR
Wissensrepräsentation und Metainformation	Ist auf reine Suche ausgelegt. Stark heterogen, Qualität schwankt, Spam, Verlinkung von Dokumenten	Auf Unternehmensprozesse ausgelegt, heterogen strukturiert, dynamisch, wenig Metainformation
Umfang der Systeme	unendlich	Eher klein (Unternehmen)
Anfragesprache	Einfach, natürliche Sprache	Advanced Search, Browsing, Kategorien
Such-Kompetenz der Anwender	Beliebig: gering -> trainiert	Gering, untrainierte Enduser
Fachkompetenz der Anwender	Beliebig: gering -> trainiert	Hoch, professionelle Anwender mit meist sehr spezifischem Fachwissen