

UNIVERSITATEA BUCUREȘTI
FACULTATEA DE MATEMATICĂ ȘI INFORMATICĂ

LUCRARE DE DISERTAȚIE

Metode de ierarhizare în Regăsirea Informației

Autor:
Călin AVASÎLCĂI

Coordonator:
Lect. Dr. Marius POPESCU

MAI 2010

Cuprins

1	Introducere	5
1.1	Scurt istoric	5
1.2	Descrierea procesului	5
1.3	Un exemplu	6
2	Regăsirea informației - generalități	7
2.1	Indexare	7
2.2	Cautare	7
2.3	Regăsirea booleană și ad-hoc	7
3	Evaluarea în regăsirea informației	9
3.1	Evaluarea unui sistem RI	9
3.2	Colecții sandard de test	10
3.3	Evaluarea rezultatelor neordonate	10
3.4	Evaluarea rezultatelor ordonate	12
4	Metode de ierarhizare	17
4.1	Modelul spațiului vectorial (VSM)	17
4.1.1	Similaritatea cosinus și td-idf	17
4.2	Euristici pentru eficientizare	17
4.3	Cluster pruning	17
4.4	Modelul axiomatic	17
4.4.1	F2-EXP	17
4.5	Modelul probabilistic	17
4.5.1	Principiul probabilistic de ierarhizare	17
4.5.2	Modelul binar de independență	18
4.5.3	Okapi BM25	21
5	Metode de agregare	23
5.1	Borda	23
5.2	Agregarea Rank Distance	23
6	Studiu comparativ	25
6.1	Unelte	25
6.1.1	Apache Lucene	25
6.1.2	Pachetul Benchmark	27
6.1.3	Apache Open Relevance Project	27
6.2	Construirea framework-ului	27
6.3	Rezultate	27
7	Concluzii	29

Capitolul 1

Introducere

Singtagma *regăsirea informației* (RI) are un spectru larg de înțelesuri. Scoaterea unui card de credit din portofel pentru a completa numărul cardului într-un formular este o formă de regăsire a informației. Dar, ca domeniu de studii academice, *regăsirea informației* poate fi definită astfel[1]:

Regăsirea informației se referă la găsirea de material (de obicei documente) de natură nestructurată (de obicei text) care satisface o nevoie de informație, din colecții mari de date (eventual stocate pe calculatoare).

Regăsirea informației obișnuia să fie o activitate practică de câțiva oameni cum ar fi bibliotecarii. Acum lumea s-a schimbat și sute de milioane de oameni folosesc sisteme de regăsire a informației precum motoare de căutare web în fiecare zi. Regăsirea informației devine rapid forma dominantă de accesare a informației întrecând modul tradițional de căutare de tip "baze de date" (căutarea după un număr de identificare).

RI este interdisciplinară, bazată pe informatică, matematică, știința informației, arhitectura informației, psihologie cognitivă, lingvistică, statistică, etc. Sisteme automate de RI au fost folosite pentru a reduce ceea ce se numește "supraîncărcarea de informație". Multe universități și librării publice folosesc sisteme RI pentru a facilita accesul la cărți, jurnale și alte tipuri de documente. Cele mai folosite aplicații RI sunt motoarele de căutare web.

1.1 Scurt istoric

Ideea de a folosi calculatoarele pentru a căuta "bucăți" de informație a fost popularizată în articolul *As We May Think* de către *Vannevar Bush* în anul 1945[2]. Primele sisteme automate de regăsire a informației au fost introduse în anii 1950 și 1960. Până în 1970 câteva tehnici au fost testate cu succes pe corpusuri mici de text cum ar fi colecția *Cranfield* (câteva mii de documente). Sisteme mari de RI cum ar fi *Lockheed Dialog system*, au fost date în folosință la începutul anilor 1970.

În 1992 Departamentul de Apărare al Statelor Unite împreună cu Institutul Național de Standarde și Tehnologie (NIST) au sponsorizat *TREC* (**T**ext **R**etrieval **C**onference) ca parte din programul TIP-STER. Scopul a fost să se asigure comunității RI infrastructura necesară pentru testarea și evaluarea metodelor de regăsire a textului pe colecții foarte mari. Acest program a acționat ca un catalizator în ceea ce privește cercetarea metodelor RI care scalează la colecții foarte mari. Apariția motoarelor de căutare pe web a adus o nevoie și mai mare de sisteme RI care pot face față unor colecții imense de date.

1.2 Descrierea procesului

Un proces de regăsire a informației începe când un utilizator introduce o interogare (*query*) în sistem. Interogările sunt formulări formale ale *nevoilor de informație*, de exemplu, cuvinte scrise în caseta de căutare a unui motor de căutare web. În RI o interogare nu identifică în mod unic un obiect din colecția de obiecte. În schimb, mai multe obiecte pot fi considerate ca răspunsuri la interogare, eventual cu grade de relevanță diferite.

Un obiect este o entitate stocată în sistemul RI. În funcție de aplicație, obiectele pot fi texte, documente, imagini, videoclipuri, etc. În general, documentele propriu-zise nu sunt ținute în sistem ci sunt reprezentate de surrogat(rezultate în urma procesării documentelor) și metadate.

Majoritatea sistemelor RI calculează un scor numeric ce reprezintă cât de bine se potrivește un obiect cu interogarea utilizatorului și apoi ierarhizează obiectele în funcție de această valoare. Cele mai bune obiecte sunt apoi afișate utilizatorului.

1.3 Un exemplu

Capitolul 2

Regăsirea informației - generalități

2.1 Indexare

2.2 Cautare

2.3 Regăsirea booleană și ad-hoc

Capitolul 3

Evaluarea în regăsirea informației

Există multe alternative pentru a proiecta un sistem RI. Cum hotărâm care dintre aceste tehnici sunt eficiente în anumite aplicații? Este indicată folosirea unui *stop list* sau a unui *stemmer*? Regăsirea informației a evoluat ca o disciplină empirică, necesitând evaluări atente și complete pentru a demonstra superioritatea performanțelor noilor tehnici apărute pe colecții reprezentative de documente.

3.1 Evaluarea unui sistem RI

Pentru a măsura eficiența unui sistem ad-hoc RI este necesară o colecție de test alcătuită din trei componente:

1. O colecție de documente
2. O serie de *nevoi de informație* exprimate prin interogări
3. Un set de judecăți de relevanță, de obicei o valoare binară (*relevant/nerelevant*) pentru fiecare pereche interogare-document.

Abordarea standard în evaluarea unui sistem RI se învârtă în jurul noțiunii de documente *relevante* și *nerelevante*. Cu privire la o nevoie de informație, unui document dintr-o colecție de test i se dă o clasificare binară: relevant sau nerelevant. Colecția de documente și setul de interogări trebuie să fie destul de mari: este nevoie de efectuarea unei medii a rezultatelor de performanță peste seturi mari de test întrucât acestea variază destul de mult de la interogare la interogare. Ca o regulă derivată din experimente, 50 de nevoi de informație reprezintă un minim suficient.

Relevanța este evaluată relativ la nevoia de informație, nu la interogare. De exemplu, o nevoie de informație ar putea fi:

Informație cu privire la faptul că vinul roșu este mai eficient în prevenirea atacurilor de cord decât vinul alb.

Interogarea asociată acestei nevoi de informație ar putea fi:

vin ȘI roșu ȘI alb ȘI atac ȘI cord ȘI eficient

Un document este relevant dacă se adresează nevoii de informație nu doar pentru că se întâmplă să conțină cuvintele din interogare. Acest lucru este des neînțeles în practică, deoarece nevoia de informație este "ascunsă". În orice caz, ea este prezentă. Pentru o interogare care constă într-un cuvânt, este foarte greu pentru un sistem RI să își dea seama care este nevoia de informație. Dar utilizatorul are o astfel de nevoie și va judeca rezultatele pe baza acesteia. Pentru a evalua un sistem este nevoie de o exprimare "deschisă" a nevoii de informație, care poate fi folosită la judecarea documentelor întoarse de sistem ca relevante sau nerelevante. Relevanța poate fi gândită ca o scală, cu unele documente *mai* relevante decât altele. Dar, pentru început, se va folosi o valoare binară a relevanței.

Multe sisteme conțin diferiți parametri de calibrare ce pot fi ajustați pentru a mări performanța. Este greșit să se raporteze rezultate pe o colecție de test obținute în urma modificării acestor parametri cu scopul de a maximiza performanța pe colecția respectivă deoarece parametrii vor fi setați să maximizeze performanța pe un anumit set de interogări în loc să facă acest lucru pentru un eșantion aleator de interogări.

3.2 Colecții standard de test

Aceasta este o listă de colecții standard de test și evaluare. Majoritatea sunt menite pentru a testa sisteme ad-hoc de RI dar sunt și câteva folosite pentru clasificare de text.

Cranfield Această colecție a fost prima care a permis măsuri precise de eficiență ale sistemelor de regăsire a informației, dar acum este prea mică pentru un experiment "serios". A fost alcătuită în anii 1950 în Regatul Unit și conține 1398 de sumare de articole despre aerodinamică, un set de 225 de interogări și judecăți de relevanță complete pentru toate perechile (interogare, document).

TREC NIST a construit un framework de evaluare pentru RI începând cu 1992. Din acest framework cele mai cunoscute sunt cele folosite la testele *TREC Ad Hoc track* în timpul primelor opt conferințe TREC (Text REtrieval Conference) între 1992 și 1999. În total, aceste colecții de test însumează 6 CD-uri și conțin 1,89 milioane de documente (compuse cu preponderență din articole de știri) și judecăți de relevanță pentru 450 de nevoi de informație numite subiecte. Colecțiile TREC 6-8 sunt formate din 150 de nevoi de informație și 528.000 de articole de știri. Acesta este probabil cea mai bună subcolecție datorită faptului că este cea mai mare și subiectele sunt mai consistente. Din cauza faptului că aceste colecții sunt foarte mari, nu există judecăți exhaustive ci sunt specificate doar pentru un subset de documente (primele k documente întoarse de sistemul pentru care nevoia de informație a fost dezvoltată).

GOV2 În ultimii ani, NIST a făcut evaluări pe colecții mult mai mari, incluzând colecția de 25 de milioane de pagini web, *GOV2*. *GOV2* este acum cea mai mare colecție de pagini web disponibilă pentru cercetare. Cu toate acestea, *GOV2* încă este de cel puțin 2 ori mai mică decât colecția de documente indexate de companiile mari de căutare web.

NTCIR Proiectul *NTCIR - NII Test Collections for IR Systems* a construit diferite colecții de test de mărimi similare cu colecțiile TREC. Aceste colecții sunt concentrate pe limba est-asiatică și pe *regăsirea informației cross-language* (interogările sunt făcute într-o limbă peste un set de documente scrise în diferite limbi).

CLEF Seria de evaluări CLEF (Cross Language Evaluation Forum) s-a concentrat pe limbile europene și regăsirea informației de tip cross-language.

Reuters Pentru clasificarea de text, cea mai folosită colecție de test a fost *Reuters-21578* compusă din 21578 articole de știri. Mai recent Reuters a publicat *Reuters Corpus Volume 1 (RCV1)*, o colecție mult mai mare constând 806.791 documente.

3.3 Evaluarea rezultatelor neordonate

Având aceste ingrediente, cum este măsurată eficiența unui sistem? Cele mai frecvente și fundamentale două măsuri în RI sunt așa-numitele *precizia* (*precision* - P) și *returnarea* (*recall* - R). Acestea sunt mai întâi definite pentru cazul în care sistemul RI întoarce un set neordonat de documente pentru o interogare și apoi extinse la liste ierarhizate de documente.

Precizia este dată de raportul dintre documentele regăsite care sunt relevante și documentele regăsite:

$$P = \frac{\#(\text{documente regăsite și relevante})}{\#(\text{documente regăsite})} = P(\text{relevante}|\text{regăsite}) \quad (3.1)$$

Returnarea este raportul dintre documentele regăsite care sunt relevante și documentele relevante:

$$R = \frac{\#(\text{documente regăsite și relevante})}{\#(\text{documente relevante})} = P(\text{regăsite}|\text{relevante}) \quad (3.2)$$

Aceste noțiuni pot fi explicate mai clar folosind tabelul de contingență de mai jos.

Atunci:

$$P = \frac{tp}{tp + fp} \quad R = \frac{tp}{tp + fn} \quad (3.3)$$

	Relevante	Nerelevante
Regăsite	pozitive adevărate(tp)	pozitive false(fp)
Neregăsite	negative false(fn)	negative adevărate(tn)

O alternativă evidentă este judecarea sistemului după acuratețe, adică, procentul clasificărilor corecte. Folosind tabelul de mai sus, $ac = (tp + tn)/(tp + fp + fn + tn)$. Acest lucru pare plauzibil din moment ce există două clase, relevant și nerelevant, iar un sistem RI poate fi privit ca un clasificator pe două clase (întoarce documentele etichetate cu *relevant*). Aceasta este, de fapt, măsura de eficiență folosită de obicei în evaluarea problemelor de clasificare.

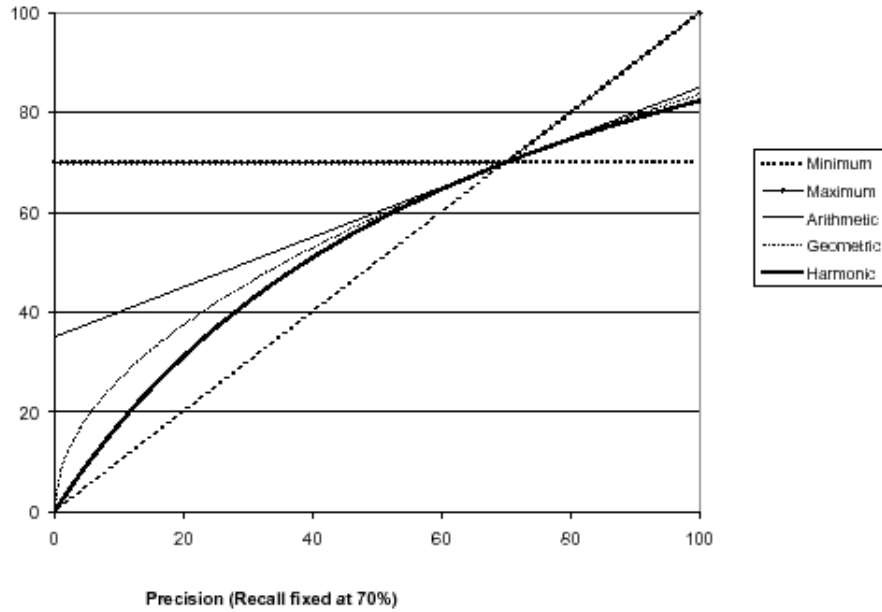


Figura 3.1: Grafic care compară media armonică cu celelalte medii

Există un motiv din cauza căruia acuratețea nu este o măsură corespunzătoare pentru regăsirea informației. În aproape toate circumstanțele, datele sunt extrem de "denaturate": de obicei, peste 99,9% din documente sunt în categoria nerelevant. Un sistem configurat să maximizeze acuratețea poate părea că funcționează foarte bine prin considerarea tuturor documentelor ca fiind nerelevante la orice întrebare. Acest lucru este inacceptabil pentru un sistem RI. Un utilizator va dori întotdeauna să vadă niste documente și poate avea o mică toleranță pentru pozitive false dacă majoritatea documentelor îi satisfac nevoia de informație. Măsurile P și R concentrează evaluarea pe întoarcerea de pozitive adevărate, întrebând ce procent de documente relevante a fost găsit și cate pozitive false au fost întoarse.

Avantajul faptului că există două măsuri diferite este acela că în cele mai multe circumstanțe una este mai importantă decât cealaltă. Un utilizator web obișnuit dorește ca fiecare rezultat de pe prima pagină să fie relevant (precizie mare) dar nu are nici cel mai mic interes să știe și mai ales să privească fiecare document relevant. Pe de altă parte, un profesionist care caută documente într-un sistem RI intern va încerca să găsească toate documentele relevante la nevoia lui de informație (returnare mare) și va tolera valori mici ale preciziei pentru a își atinge scopul. În orice caz, există un compromis între cele două valori: se poate întotdeauna obține un recall maxim dar precizie mică dacă se întorce o listă cu toate documentele. Recall-ul este o funcție non-descrescătoare de numărul de documente regăsite. Pe de altă parte, într-un sistem bun, precizia de obicei scade odată cu creșterea numărului de documente regăsite. În general se dorește o anumită valoare a recall-ului, tolerându-se un anumit procent de pozitive false.

O măsură care înglobează atât precizia cât și returnarea este *măsura F*, care se calculează ca media armonică ponderată a celor două valori:

$$F = \frac{1}{\alpha \frac{1}{P} + (1 - \alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R} \text{ unde } \beta^2 = \frac{1 - \alpha}{\alpha} \quad (3.4)$$

unde $\alpha \in [0, 1]$ și, ca urmare, $\beta^2 \in [0, \infty)$. Măsura F balansată ținând cont în mod egal de precizie și $\text{recall}(\alpha = 1/2, \beta = 1)$. Este notată de obicei cu F_1 (de la $F_{\beta=1}$):

$$F_{\beta=1} = \frac{2PR}{P + R} \quad (3.5)$$

Valori ale lui β mai mici decât 1 scot în evidență precizia, în timp ce valori mai mari ca 1 pun accentul pe returnare. Valorile discutate până acum sunt măsuri între 0 și 1, dar se mai scriu și ca procente pe o scală de la 1 la 100. Media armonică este folosită în locul mediei aritmetice deoarece în cazul în care un sistem ar întoarce toate documentele (100% recall), ar avea tot timpul cel puțin $F = 50\%$. Media armonică a două numere ($a < b$) este mai apropiată de minim decât de media lor aritmetică (vezi figura 3.1).

3.4 Evaluarea rezultatelor ordonate

Măsurile *precision*, *recall* și F sunt bazate pe seturi; sunt calculate folosind seturi neordonate de documente. Este nevoie de o extindere a acestor măsuri (sau de o definiție de noi măsuri) dacă se dorește evaluarea unor rezultate ierarhizate care reprezintă standardul în zilele noastre. Într-un context de regăsire ierarhizată setul care interesează este dat de *primele k* documente. Pentru fiecare astfel de set, valorile preciziei și returnării pot fi reprezentate grafic printr-o *curbă precision-recall*, ca cea din figura 3.2.

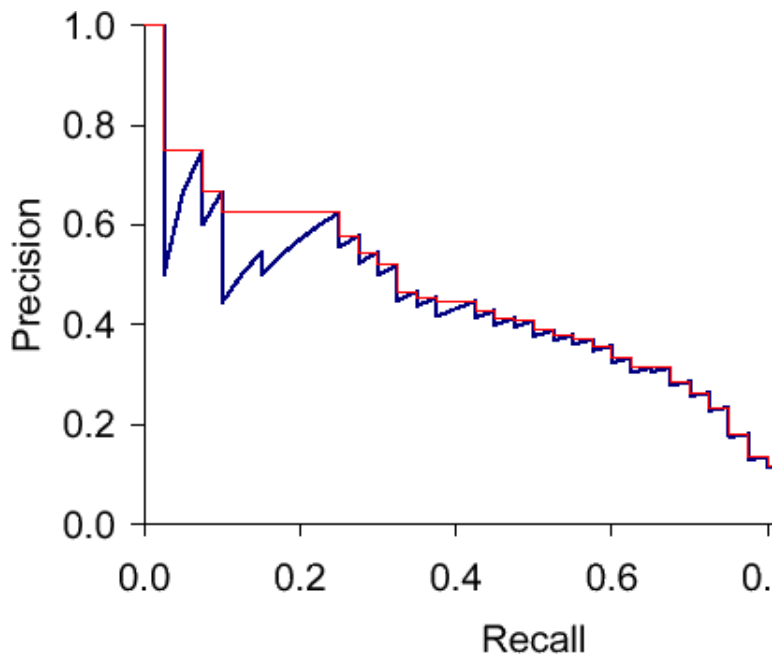


Figura 3.2: Grafic *precision-recall*.

Curbele *precision-recall* au o formă de "fierăstrău": dacă al $(k + 1)$ -lea document regăsit este nerlevant, atunci *recall*-ul este același ca pentru primele k documente regăsite, dar *precizia* scade. Dacă documentul este relevant atunci ambele valori cresc. Este de obicei util să se elimine această formă și modul standard de a face acest lucru este folosind *precizia interpolată* p_{interp} . La un anumit nivel de *recall* r , p_{interp} este dat de cea mai mare valoare a preciziei găsită pentru orice nivel de *recall* $r' \geq r$:

$$p_{interp} = \max_{r' \geq r} p(r') \quad (3.6)$$

Precizia interpolată este ilustrată în figura 3.2.

Examinarea întregii curbe *precizie-recall* este foarte informativă, dar, câteodată este necesar ca informația să fie redusă la câteva numere. Modul tradițional de a realiza acest lucru (folosit, de exemplu,

la primele 8 evaluari TREC Ad Hoc) este *precizia medie interpolată în 11 puncte*. Pentru fiecare nevoie de informație, precizia interpolată este măsurată în cele 11 puncte de recall: 0.1, 0.2, ..., 1.0. Exemplul din figură 3.2 este ilustrat în tabelul 3.1.

Recall	Precizie interpolată
0.0	1.00
0.1	0.67
0.2	0.63
0.3	0.55
0.4	0.45
0.5	0.41
0.6	0.36
0.7	0.29
0.8	0.13
0.9	0.10
1.0	0.08

Tabelul 3.1: Calcularea preciziei medii interpolate în 11 puncte.

Pentru fiecare nivel de recall se calculează media aritmetică a preciziei interpolate la acel nivel de recall pentru fiecare nevoie de informație din colecție. O curbă precizie-recall cu 11 puncte poate fi vizualizată apoi grafic. Un exemplu este prezent în figura 3.3.

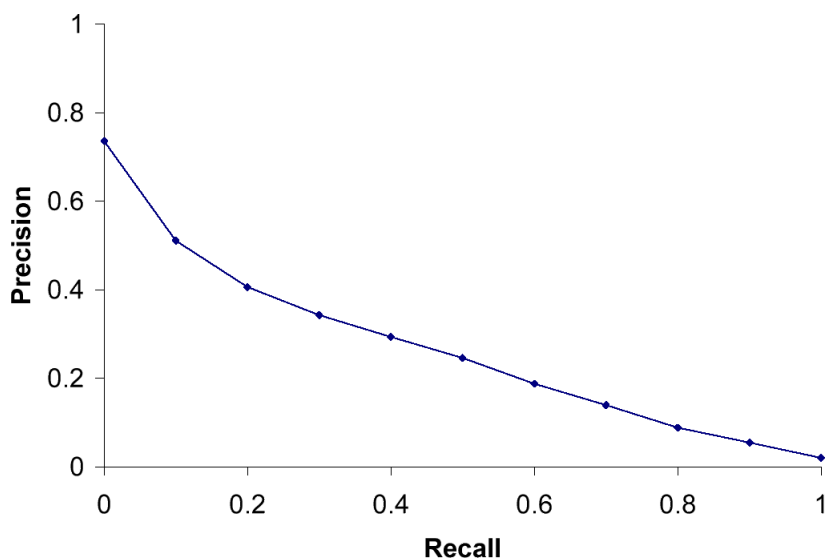


Figura 3.3: Grafic al preciziei medii interpolate în 11 puncte.

În ultimii ani, alte măsuri au devenit mai comune. Cea mai folosită în comunitatea TREC este *Mean Average Precision (MAP)*, care pune la dispoziție o singură măsură de calitate de-a lungul nivelelor de recall. Printre tipurile de masuri de evaluare s-a arătat că MAP funcționează foarte bine la capitolele discriminare și stabilitate. Pentru o singură nevoie de informație, precizia medie este media valorilor de precizie obținute pentru primele k documente după regăsirea fiecărui document relevant. Media acestei valori peste toate nevoile de informație este MAP. Dacă setul de documente relevante pentru o nevoie de informație $q_j \in Q$ este d_1, \dots, d_{m_j} și R_{jk} este setul de rezultate de la primul până la documentul d_k , atunci:

$$MAP(Q) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{m_j} \sum_{k=1}^{m_j} Prec(R_{jk}) \quad (3.7)$$

Când un document relevant nu este regăsit, precizia este considerată 0. Pentru o singură nevoie de informație, precizia medie aproximează aria de sub curba neinterpolată precizie-recall, ășadar MAP aproximează aria medie de sub curba precizie-recall pentru un set de interogări.

Folosind MAP, nu sunt alese nivele fixe de recall și nu este nevoie de interpolare. Un set de nevoi de informație trebuie să fie cuprinzător și diversificat pentru a putea măsura cât mai exact eficiența unui sistem.

Măsurile descrise până au în componența precizia pentru fiecare nivel de recall. În cazul multor aplicații, cum ar fi un motor de căutare web, acest lucru nu reflectă neapărat nevoile unui utilizator. În acest caz contează mai mult câte rezultate bune se găsesc în primele pagini. Acest lucru conduce la nevoia de a măsura precizia pentru primele documente regăsite (10 sau 30). Această tehnică poartă numele de "*Precizia la k*". Are avantajul că nu necesită o estimare a dimensiunii setului de documente relevante (m_j în cazul MAP) și dezavantajele că este cea mai puțin stabilă dintre metodele de evaluare și că nu se poate calcula media prea bine pentru ea.

O alternativă care atenuează această problemă este *R-precizia*. Implică existența unui set *Rel* de documente despre care se știe că sunt relevante, din care se calculează apoi precizia pentru primele $|Rel|$ documente regăsite. R-precizia se ajustează la mărimea setului de documente relevante: un sistem perfect ar putea puncta 1 la această măsură pentru fiecare interogare, în timp ce până și un sistem perfect, ar putea să atingă o precizie la 20 de 0.4 dacă ar exista doar 8 documente în colecția de documente relevante. Calcularea mediei peste setul de interogări are mai mult sens în acest caz. Dacă există $|Rel|$ documente relevante pentru o interogare, se examinează primele $|Rel|$ rezultate date de un sistem și se găsesc r documente relevante, atunci precizia (și deci R-precizia) este $r/|Rel|$, dar și recall-ul este tot $r/|Rel|$. Așadar, R-precizia este identică cu o altă măsură folosită câteodată: *break-even point*, măsură definită ca valoarea la care precizia și recall-ul sunt egale. Ca și *Precizia la k*, *R-precizia* descrie un singur punct de pe curbă, și e uneori neclar de ce interesează mai mult nivelul în care cele două valori sunt egale decât nivelul cel mai bun de pe curbă (în care F este maxim) sau un nivel de interes pentru o anumită aplicație (*Precizia la k*). Cu toate acestea, R-precizia pare să fie corelată cu MAP, lucru constatat în urma experimentelor.

Un alt concept folosit uneori în evaluarea unui sistem este *curba ROC* ("*ROC*" vine de la *Receiver Operating Characteristics*). O curbă ROC reprezintă grafic rata pozitivelor adevărate (*sensitivitate*) ca funcție de rata pozitivelor false ($1 - \text{specificitate}$). Aici sensibilitate este un alt termen pentru recall. Rata de pozitive false este dată de $fp/(fp + tn)$. Figura 3.4 ilustrează curba ROC corespunzătoare cubei precizie-recall din figura 3.2.

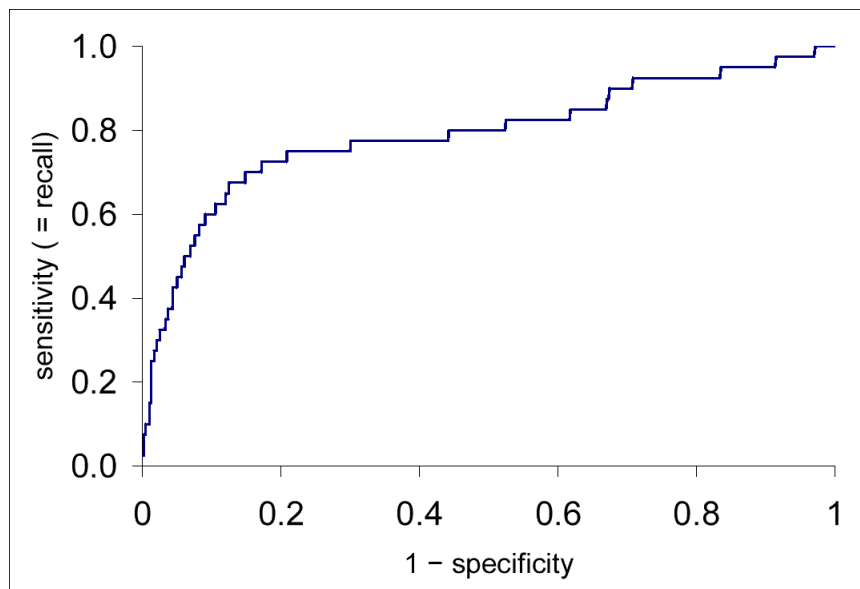


Figura 3.4: Curba ROC

O curbă ROC urcă întotdeauna din stânga jos spre dreapta sus. Pentru un sistem bun, graficul urcă abrupt în partea stângă. Are mai mult sens când se ia în considerare întreg spectrul de regăsire și pune la dispoziție o altă perspectivă asupra datelor. În unele cazuri se folosește estimarea ariei de sub curba ROC, care este o valoare analoagă măsurii MAP.

O altă abordare care a fost adoptată din ce în ce mai des, în special în cazul sistemelor care folosesc

machine learning, este măsura *câștigului cumulativ* (*cumulative gain*) și, în particular, *normalized discounted cumulative gain* (*NDCG*). NDCG este proiectat pentru situații în care relevanța este nebinară. Ca și *Precizia la k*, este evaluată peste primele k rezultate. Pentru un set de interogări Q, fie $R(j, d)$ scorul de relevanță pe care evaluatorii l-au dat documentului d pentru interogarea j. Atunci,

$$NDCG(Q, k) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} Z_{kj} \sum_{m=1}^k \frac{2^{R(j,m)} - 1}{\log_2(1 + m)}, \quad (3.8)$$

unde Z_{kj} este un factor de normalizare calculat pentru a face ca NDCG-ul unei ierarhizări perfecte la k pentru interogarea j să fie 1. Pentru interogări pentru care $k' < k$ documente sunt returnate, ultima sumă se face până la k'

Capitolul 4

Metode de ierarhizare

4.1 Modelul spațiului vectorial (VSM)

4.1.1 Similaritatea cosinus și td-idf

4.2 Euristici pentru eficientizare

4.3 Cluster pruning

4.4 Modelul axiomatic

4.4.1 F2-EXP

4.5 Modelul probabilistic

Dacă s-ar cunoaște relevanța unui subset de documente, s-ar putea estima probabilitatea apariției unui termen t într-un document relevant $P(t|R = 1)$ și, ca urmare, acesta ar putea reprezenta baza unui clasificator care decide dacă un document este relevant sau nu.

Utilizatorii încep cu *nevoi de informație* pe care le transformă în *interogări*. În mod similar, documentele sunt transformate în *reprezentări de documente* care diferă de primele cel puțin prin felul în care textul este împărțit în token-i. Bazându-se pe aceste două reprezentări, un sistem încearcă să determine cât de bine satisfac documentele nevoile de informații. În modelul Boolean sau VSM, dându-se numai o interogare, pentru un sistem RI nevoia de informație este incertă. Dându-se interogarea și reprezentarea documentelor, un sistem trebuie să "ghicească" dacă un document are conținut relevant pentru respectiva nevoie de informație. Teoria probabilităților pune la dispoziție o fundație de principii potrivite pentru raționament în situații incerte. Aceste principii pot fi exploatate pentru a estima cât de probabil este ca un document să fie relevant pentru o nevoie de informație.

Există mai multe posibile modele probabilistice de regăsire. În continuare voi discuta despre *principiul probabilistic de ierarhizare* și despre *modelul binar de independență*, care a fost primul model probabilistic de regăsire. În final voi prezenta și sistemul de ponderare *Okapi BM25*, care a avut un succes destul de mare în practică.

În acest context, este util conceptul de șanse (*odds*), pe lângă cel de probabilitate.

$$\text{Odds} : O(A) = \frac{P(A)}{P(\bar{A})} = \frac{P(A)}{1 - P(A)}. \quad (4.1)$$

4.5.1 Principiul probabilistic de ierarhizare

Cazul 1/0 loss

Presupunem că sistemul de RI întoarce ca răspuns la o interogare o listă ordonată de documente și folosirea unei notații binare pentru relevanță. Pentru o interogare q și un document d , fie $R_{d,q}$ o vari-

abilă aleatoare care indică dacă d este relevant în contextul interogării q . Variabila ia valoarea 1 când documentul este relevant și 0 altfel.

Folosind un model probabilistic, ordinea evidentă în care documentele trebuie prezentate utilizatorului este dată de ierarhizarea documentelor după probabilitatea estimată de relevanță în raport cu nevoia de informație: $P(R_{d,q} = 1|d, q)$. Vom scrie R în loc de $R_{d,q}$. Aceasta reprezintă temelia *principiului probabilistic de ierarhizare (PRP)*:

Dacă răspunsul unui sistem la fiecare cerere este o ierarhizare a documentelor din colecție în ordinea descrescătoare a probabilității de relevanță, unde probabilitățile sunt estimate cât mai bine cu putință pe baza datelor pe care sistemul le are la îndemână, eficiența sistemului este cea mai bună care se poate obține folosind aceste date.

În cel mai simplu caz al PRP, nu există costuri de regăsire sau alte motive de îngrijorare care să valorifice diferit acțiunile sau erorile. Se pierde un punct fie pentru întoarcerea unui document nerelevant, fie pentru lipsa întoarcerii unui document relevant. O astfel de evaluare binară asupra preciziei poartă numele de *1/0 loss*. Scopul este să se întoarcă cele mai bune k rezultate posibile, pentru orice valoare k aleasă de utilizator. PRP spune că documentele trebuie ierarhizate în ordinea descrescătoare a $P(R = 1|d, q)$.

Teoremă 1 *PRP este optim în sensul că minimizează pierderea așteptată (sau riscul Bayes) în cazul 1/0 loss.*

Această teoremă este adevărată dacă toate probabilitățile sunt corecte, ceea ce în practică este imposibil. Cu toate acestea, PRP reprezintă o fundație pentru contruirea de modele de RI.

Costuri de regăsire

Să presupunem existența unui model de costuri de regăsire. Fie C_1 costul de regăsire a unui document relevant și C_0 costul de regăsire a unui document nerelevant. Atunci pentru un document d și pentru toate documentele d' neregăsite dacă

$$C_1 \times P(R = 1|d) + C_0 \times P(R = 0|d) \leq C_1 \times P(R = 1|d') + C_0 \times P(R = 0|d') \quad (4.2)$$

atunci d este următorul document care trebuie întors. Acest model asigură un cadru formal în care putem modela costurile diferențiale ale falselor-pozitive și falselor-negative.

4.5.2 Modelul binar de independență

BIM este modelul care a fost folosit cu PRP. Introduce câteva asumții simple care permit estimarea funcției probabilistice $P(R|d, q)$. Aici binar este echivalent cu boolean: atât documentele cât și interogările sunt reprezentate ca vectori binari de incidență a termenilor. Un document d este reprezentat de vectorul $\vec{x} = (x_1, \dots, x_M)$, unde $x_t = 1$ dacă termenul t este prezent în documentul d și $x_t = 0$ altfel. În contextul acestei reprezentări, multe documente pot avea aceeași reprezentare. În mod similar, interogarea q este reprezentată prin vectorul de incidență \vec{q} . "Independență" se referă la faptul că termenii sunt modelați așa cum apar în documente în mod independent. Modelul nu recunoaște niciun tip de asociere între termeni. Această asumție este, evident, incorectă, dar, în pofida acestui aspect, oferă rezultate satisfăcătoare în practică. Este asumția care stă și la baza modelului *Bayes Naiv*, și este într-un fel echivalentă cu asumția din VSM, în care fiecare termen reprezintă o dimensiune ortogonală față de celelalte.

Pentru a face o strategie probabilistică de regăsire precisă, trebuie estimat modul în care termenii din documente contribuie la relevanță. Cu alte cuvinte, trebuie să specificăm cum frecvența termenilor într-un document, numărul de documente care conțin un termen, lungimea unui document și alte statistici influențează calculul relevanței unui document. După acest proces documentele vor fi ordonate în ordinea descrescătoare a acestor probabilități estimate.

Se pornește de la asumția că relevanța fiecărui document este independentă de relevanța celorlalte documente. În practică, acest lucru pune o problemă în momentul în care sunt întoarse documente duplicate sau aproape duplicate. În contextul BIM, probabilitatea că un document este relevant la o

interogare $P(R|d, q)$ este modelată via probabilitatea $P(R|\vec{x}, \vec{q})$, folosind vectorii de incidență. Apoi, aplicând regula lui Bayes, se obține:

$$P(R = 1|\vec{x}, \vec{q}) = \frac{P(\vec{x}|R = 1, \vec{q})P(R = 1|\vec{q})}{P(\vec{x}, \vec{q})} \quad P(R = 0|\vec{x}, \vec{q}) = \frac{P(\vec{x}|R = 0, \vec{q})P(R = 0|\vec{q})}{P(\vec{x}, \vec{q})} \quad (4.3)$$

Aici, $P(\vec{x}|R = 1, \vec{q})$ și $P(\vec{x}|R = 0, \vec{q})$ reprezintă probabilitatea ca dacă un document relevant, respectiv nerelevant, este întors, acesta să aiba reprezentarea \vec{x} . Aceste probabilități nu se pot calcula exact, așa că trebuie folosiți estimatori: statistici ale colecției de documente sunt folosite pentru a estima aceste probabilități. $P(R = 1|\vec{q})$ și $P(R = 0|\vec{q})$ reprezintă probabilitatea apriori de a întoarce un document relevant, respectiv nerelevant, dată fiind interogarea \vec{q} . Pentru că un document este fie relevant fie nerelevant în contextul unei interogări, avem:

$$P(R = 1|\vec{x}, \vec{q}) + P(R = 0|\vec{x}, \vec{q}) = 1. \quad (4.4)$$

Derivarea unei funcții de ierarhizare

Decât să estimăm $P(R = 1|\vec{x}, \vec{q})$ direct, deoarece interesează doar ordinea în care sunt întoarse documentele, folosim alte cantități care sunt mai ușor de calculat și care au ca rezultat aceeași ordine. Putem ordona documentele după șansele (odds) de relevanță, ceea ce duce la o simplificare a relației:

$$O(R|\vec{x}, \vec{q}) = \frac{P(R = 1|\vec{x}, \vec{q})}{P(R = 0|\vec{x}, \vec{q})} = \frac{\frac{P(\vec{x}|R=1, \vec{q})P(R=1|\vec{q})}{P(\vec{x}, \vec{q})}}{\frac{P(\vec{x}|R=0, \vec{q})P(R=0|\vec{q})}{P(\vec{x}, \vec{q})}} = \frac{P(R = 1|\vec{q})}{P(R = 0|\vec{q})} \cdot \frac{P(\vec{x}|R = 1, \vec{q})}{P(\vec{x}|R = 0, \vec{q})} \quad (4.5)$$

Termenul stâng al expresiei din dreapta al ecuației 4.5 este constant pentru o interogare dată. Pentru că interesează doar ordinea, nu e nevoie să se estimeze. Trebuie estimat, în schimb, celălalt termen, lucru care pare dificil inițial: cum poate fi precis estimată probabilitatea unui întreg vector de incidență? Pentru a face posibilă estimarea, se face asumția de *condițional-independență Naive Bayes*: prezența sau absența unui cuvânt într-un document este independentă de prezența sau absența oricărui alt cuvânt:

$$\frac{P(\vec{x}|R = 1, \vec{q})}{P(\vec{x}|R = 0, \vec{q})} = \prod_{t=1}^M \frac{P(x_t|R = 1, \vec{q})}{P(x_t|R = 0, \vec{q})} \quad (4.6)$$

Așadar:

$$O(R|\vec{x}, \vec{q}) = O(R|\vec{q}) \cdot \prod_{t=1}^M \frac{P(x_t|R = 1, \vec{q})}{P(x_t|R = 0, \vec{q})}. \quad (4.7)$$

Pentru că fiecare x_t este fie 0, fie 1, putem separa termenii astfel:

$$O(R|\vec{x}, \vec{q}) = O(R|\vec{q}) \cdot \prod_{t:x_t=1} \frac{P(x_t = 1|R = 1, \vec{q})}{P(x_t = 1|R = 0, \vec{q})} \cdot \prod_{t:x_t=0} \frac{P(x_t = 0|R = 1, \vec{q})}{P(x_t = 0|R = 0, \vec{q})}. \quad (4.8)$$

Fie $p_t = P(x_t = 1|R = 1, \vec{q})$ probabilitatea ca termenul x_t să apară într-un document relevant la interogarea dată și $u_t = P(x_t = 1|R = 0, \vec{q})$, probabilitatea ca termenul să apară într-un document nerelevant. Aceste cantități pot fi vizualizate în tabelul de contingență următor (suma pe coloane este 1):

document		relevant (R=1)	nerelevant(R=0)
termen prezent	$x_t = 1$	p_t	u_t
termen absent	$x_t = 0$	$1 - p_t$	$1 - u_t$

Facând asumția că termenii care nu apar în interogare au probabilități egale de apariție într-un document relevant, respectiv nerelevant: $q_t = 0 \Rightarrow p_t = u_t$, vor trebui luați în considerare doar termenii care apar în interogare:

$$O(R|\vec{x}, \vec{q}) = O(R|\vec{q}) \cdot \prod_{t:x_t=q_t=1} \frac{p_t}{u_t} \cdot \prod_{t:x_t=0, q_t=1} \frac{1 - p_t}{1 - u_t}. \quad (4.9)$$

Primul produs este peste termenii interogării care apar în document și produsul din dreapta peste cei care nu apar.

Expresia poate fi manipulată prin includerea termenilor găsiți în document în produsul din dreapta, dar, în același timp, ajustând produsul stâng pentru simplificare:

$$O(R|\vec{x}, \vec{q}) = O(R|\vec{q}) \cdot \prod_{t:x_t=q_t=1} \frac{p_t(1-u_t)}{u_t(1-p_t)} \cdot \prod_{t:q_t=1} \frac{1-p_t}{1-u_t}. \quad (4.10)$$

Produsul drept este acum peste toți termenii interogării, ceea ce înseamnă că e constant pentru o interogare, la fel ca $O(R|\vec{q})$. Așa dar, singura cantitate care trebuie estimată pentru a ierarhiza documentele este cea din produsul stâng. Putem ordona documentele și după rezultatul logaritmului produsului, deoarece log este o funcție monotonă. Cantitatea folosită la ierarhizare se numește *valoarea statusului de regăsire* (*RSV - retrieval status value*):

$$RSV_d = \log \prod_{t:x_t=q_t=1} \frac{p_t(1-u_t)}{u_t(1-p_t)} = \sum_{t:x_t=q_t=1} \log \frac{p_t(1-u_t)}{u_t(1-p_t)}. \quad (4.11)$$

Totul se reduce la clacularea RSV. Definim c_t :

$$c_t = \log \frac{p_t(1-u_t)}{u_t(1-p_t)} = \log \frac{p_t}{1-p_t} + \log \frac{1-u_t}{u_t}. \quad (4.12)$$

Termenii c_t reprezintă proporțiile *log odds* pentru termenii interogării. Valolare va fi 0 dacă un termen are șanse egale să apară într-un document relevant, respectiv nerelevant, și pozitivă dacă este mai probabil să apară într-un document relevant. Cantitățile c_t funcționează ca ponderi ale termenilor în model, iar scorul pentru o interogare și un document este: $RSV_d = \sum_{x_t=q_t=1} c_t$. Problema rămasă este cum să se estimeze cantitățile c_t pentru o colecție de documente și o interogare.

Estimări teoretice

Următorul tabel de contingență prezintă o serie de statistici ale colecției, unde df_t reprezintă numărul de documente care conțin termenul t :

	document	relevante	nerelevante	total
termen prezent	$x_t = 1$	s	$df_t - s$	df_t
termen absent	$x_t = 0$	$S - s$	$(N - df_t) - (S - s)$	$N - df_t$
	total	S	$N - S$	N

Așadar, $p_t = s/S$ și $u_t = (df_t - s)/(N - S)$ și

$$c_t = K(N, df_t, S, s) = \log \frac{s/(S-s)}{(df_t - s)/((N - df_t) - (S - s))}. \quad (4.13)$$

Pentru a evita posibilitatea apariției de zerouri (de exemplu, toate sau niciun document relevant conține un anumit termen) se adaugă $\frac{1}{2}$ la fiecare dintre cele 4 cantități și apoi se ajustează totalurile ($N + 2$). Ca urmare avem:

$$\hat{c}_t = K(N, df_t, S, s) = \log \frac{(s + \frac{1}{2})/(S - s + \frac{1}{2})}{(df_t - s + \frac{1}{2})/(N - df_t - S + s + \frac{1}{2})}. \quad (4.14)$$

Adgarea valorii $\frac{1}{2}$ este o formă simplă de *uniformizare*.

Estimări practice

Sub asumția că documentele relevante reprezintă un procent infim din colecție, este plauzibilă aproximarea statisticilor pentru documentele nerelevante cu statisticile pe întreaga colecție. Ca urmare, u_t (probabilitatea ca un document nerelevant să conțină termenul t pentru o interogare) poate fi aproximat cu df_t/N și:

$$\log \frac{1-u_t}{u_t} = \log \frac{N-df_t}{df_t} \approx \log \frac{N}{df_t} \quad (4.15)$$

Rezultatul este interesant și prin faptul că furnizează o justificare teoretică a celei mai întâlnite forme de ponderare *idf* folosită în VSM.

Tehnica de aproximare din ecuația 4.15 nu poate fi ușor extinsă la documente relevante. Cantitatea p_t poate fi estimată în mai multe moduri:

1. Se poate folosi frecvența termenilor din documentele cunoscute deja ca fiind relevante (dacă se cunosc - metodă folosită în cadrul *feedback-ului de relevanță*)
2. Se poate presupune că fiecare termen are șanse egale să apară într-un document relevant: $p_t = 0.5$. Această estimare este destul de slabă. Combinând această metodă cu aproximarea lui u_t de mai sus, ierarhizarea documentelor este dată de termenii interogării care apar în documente scalați cu ponderea *idf*.
3. O altă aproximare propusă folosește statisticile aparițiilor termenilor în colecție: $p_t = df_t/N$.

4.5.3 Okapi BM25

Metodele probabiliste sunt unele din cele mai vechi modele formale în RI. Încă din 1970 erau privite ca o oportunitate pentru a pune bazele teoretice în RI și, odată cu "renașterea" modelelor probabiliste în lingvistica computațională în anii 1990, această oportunitate s-a întors iar metodele probabiliste reprezintă unul dintre subiectele cele mai discutate subiecte în materie de RI. Obținerea unor aproximări rezonabile ale probabilităților necesare pentru un model RI probabilistic este posibilă, dar necesită prezumții majore. În modelul BIM acestea sunt:

- o reprezentare booleană a documentelor, interogărilor, relevanței
- independența termenilor
- termenii care nu apar în interogare nu afectează rezultatul
- valorile de relevanță ale documentelor sunt independente

Poate că din cauza severității asumpțiilor de modelare este dificilă obținerea unei performanțe mai bune. O problemă generală pare să fie că modelele probabiliste fie necesită informații parțiale de relevanță, fie duc la derivarea unor scheme aparent inferioare de ponderare a termenilor.

Această situație s-a schimbat în anii 1990 când schema de ponderare *BM25* a avut rezultate foarte bune și a început să fie adoptată de multe sisteme RI. Diferența dintre sistemele RI bazate pe *spatiul vectorial* și cele bazate pe modelul probabilistic nu este așa de mare; în ambele cazuri se construiește un sistem similar, singura diferență fiind că scorul documentelor în contextul unei interogări este dat pe de-o parte de *similaritatea cosinus* aplicată pe vectori de ponderi *tf-idf*, iar pe de altă parte de o formulă ușor diferită motivată de teoria probabilităților.

Un model nebinar

Modelul BIM a fost inițial proiectat pentru scurte înregistrări de cataloage și a funcționat destul de bine în acest context, dar pentru căutări *full-text* pe colecții mari este evident că un model trebuie să ia în considerare frecvența termenilor și lungimea documentelor. Schema de ponderare numită Okapi după sistemul în care a fost inițial implementată, a fost proiectată folosind un model probabilistic sensibil la aceste tipuri de informație fără să introducă prea mulți parametri adiționali.

Cel mai simplu scor pentru un document d este dat de adunarea ponderilor *idf* ale termenilor interogării prezenți în document:

$$RSV_d = \sum_{t \in q} \log \frac{N}{df_t} \quad (4.16)$$

Pornind de la formula din ecuația 4.14, și estimând $S = s = 0$ în absența feedback-ului de relevanță, se obține o formulare alternativă a *idf*:

$$RSV_d = \sum_{t \in q} \log \frac{N - df_t + \frac{1}{2}}{df_t + \frac{1}{2}} \quad (4.17)$$

Această variantă are un comportament ciudat: dacă un termen apare în peste jumătate din documentele din colecție, modelul dă o pondere negativă, lucru care este nedorit. În cazul folosirii unui *stop list* acest lucru nu se întâmplă de obicei.

Ecuția 4.16 poate fi îmbunătățită prin folosirea frecvenței termenilor și a lungimii documentului:

$$RSV_d = \sum_{t \in q} \log \left[\frac{N}{df_t} \right] \cdot \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b \times (L_d/L_{ave})) + tf_{td}} \quad (4.18)$$

Aici, tf_{td} este frecvența termenului t în documentul d și L_d și L_{ave} sunt lungimea documentului d , respectiv lungimea media a unui document din colecție. Variabila k_1 este un parametru pozitiv de ajustare care calibrează scalarea frecvenței tf_{td} : $k_1 = 0$ corespunde modelului binar, iar o valoare mare corespunde folosirii frecvenței brute. b este un alt parametru de calibrare ($0 \leq b \leq 1$) care determină scalarea după lungimea documentului: $b = 1$ presupune scalarea completă a ponderii termenului cu lungimea documentului în timp ce $b = 0$ presupune lipsa normalizării cu lungimea.

Dacă interogarea este lungă, atunci am putea folosi o ponderare similară pentru termenii din interogare. Acest lucru este adecvat dacă interogările au lungimi de dimensiunile unui paragraf, alfel este necesar:

$$RSV_d = \sum_{t \in q} \log \left[\frac{N}{df_t} \right] \cdot \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b \times (L_d/L_{ave})) + tf_{td}} \cdot \frac{(k_3 + 1)tf_{tq}}{k_3 + tf_{tq}}. \quad (4.19)$$

Aici tf_{tq} este frecvența termenului t în interogarea și k_3 este un alt parametru pozitiv de calibrare care ajustează scalarea frecvenței termenilor din interogare. Parametru b pentru normalizarea lungimii interogării este necesar.

Acești parametri în mod ideal sunt setați să optimizeze performanța pe o colecție de test. Căutarea valorilor care să maximizeze performanța poate fi făcută manual sau automat. În absența unor astfel de optimizări, experimentele au arătat că valori bune pentru acești parametri sunt $b = 0.75$ și $1.2 \leq k_1, k_3 \leq 2$.

Formulele BM25 de ponderare a termenilor au fost folosite cu succes pe o varietate de colecții și tipuri de căutare. Au avut o performanță extrem de bună la evaluările TREC și au fost implementate în multe sisteme RI.

Capitolul 5

Metode de agregare

5.1 Borda

5.2 Agregarea Rank Distance

Capitolul 6

Studiu comparativ

6.1 Unelte

6.1.1 Apache Lucene

Similaritate și ierarhizare

Lucene combină modelul boolean (BM) cu modelul de spațiu vectorial (VSM): documentele care trec de BM sunt etichetate cu un scor de către VSM.

În VSM, documentele și interogările sunt reprezentate ca vectori de ponderi într-un spațiu multidimensional, unde fiecare termen din index este o dimensiune și ponderile sunt valorile *tf-idf*. VSM nu necesită faptul ca ponderile să fie valori *tf-idf*, dar aceste ponderi au rezultate foarte bune în practică, și, ca urmare, *Lucene* folosește această abordare. Pentru un termen t și un document (sau interogare) x , $tf(t,x)$ crește odată cu numărul de ocurențe ale lui t în x iar $idf(t)$ descrește odată cu creșterea numărului de documente din index care îl conțin pe t .

Scorul documentului d pentru interogarea q este dat de *similaritatea cosinus* pentru vectorii de ponderi $V(q)$ și $V(d)$:

$$\cos - sim(q, d) = \frac{V(q)V(d)}{|V(q)||V(d)|},$$

unde numărătorul reprezintă produsul scalar, iar numitorul, produsul normelor euclidiene. Ecuația poate fi văzută și ca produsul scalar dintre cei doi vectori normalizați.

Lucene perfecționează VSM atât în materie de calitate cât și de uzabilitate.

- Normalizarea lui $V(d)$ la vectorul unitate poate pune unele probleme în sensul că îndepărtează toată informația despre lungimea documentului. Pentru unele documente, lucrul acesta poate reprezenta o problemă. Pentru a evita această problemă, *Lucene* folosește un alt factor de normalizare a lungimii documentului, care normalizează vectorul la un vector mai mare sau egal decât vectorul unitate: *doc-len-norm*(d).
- La indexare utilizatorii pot specifica faptul că unele documente sunt mai importante decât altele prin asignarea unui *boost* respectivelor documente. Ca urmare, scorul fiecărui document este multiplicat cu această valoare: *doc-boost*(d).
- *Lucene* este bazat pe câmpuri (secțiuni ale unui document), și, ca urmare, fiecare termen al unei interogări se aplică unui singur câmp, normalizarea vectorului se aplică la nivel de câmp, și se pot specifica și nivele de boost pentru câmpuri.
- Același câmp poate fi adăugat unui document în timpul indexării de mai multe ori, iar, ca urmare, nivelul de boost al acelui câmp este dat de înmulțirea nivelelor de boost ale adăugărilor.
- La căutare utilizatorii pot specifica nivele de boost pentru fiecare interogare, sub-interogare și termen al unei interogări.

- Un document poate fi relevant la o interogare cu mai mulți termeni fără să conțină toți termenii prezenți în interogare, iar documentele în care apar mai mulți termeni pot fi "răsplătite" printr-un factor de coordonare, care este mai mare când mai mulți termeni sunt prezenți: $\text{coord-factor}(q, d)$.

Facând asumptia simplificatoare că există un singur câmp în index, *formula conceptuala de scor* pentru Lucene este următoarea:

$$\text{score}(q, d) = \text{coordfactor}(q, d) \times \text{queryboost}(q) \times \frac{V(q) \times V(d)}{|V(q)|} \times \text{doclennorm}(d) \times \text{docboost}(d)$$

Din această formulă se derivează *formula practică de scor* care este implementată de Lucene. Pentru calcularea eficientă a scorului, unele componente sunt calculate și agregate la indexare:

- Nivelul de boost pentru interogare este cunoscut când căutarea începe.
- Norma euclidiană a vectorului interogare poate fi calculată când începe căutarea, dat fiind faptul că e independentă de documentul pentru care se calculează scorul la un moment dat. Din perspectiva optimizării, merită pusă întrebarea: *are rost să se normalizeze vectorul interogării, din moment ce toate scorurile vor fi înmulțite cu aceeași valoare?* Ca urmare, ierarhia documentelor pentru o interogare dată nu va fi afectată de normalizare. Există două motive pentru a păstra normalizarea:
 - scorurile unui document pentru interogări distincte trebuie să fie comparabile (într-o anumită măsură)
 - aplicarea normalizării păstrează scorurile "în jurul" vectorului unitate, împiedicând astfel alterarea scorurilor din cauza limitării de precizie ale numerelor în virgulă mobilă
- Norma pentru fiecare document $\text{doc-len-norm}(d)$ și nivelul de boost $\text{doc-boost}(d)$ sunt cunoscute la indexare. Sunt calculate și rezultatul înmulțirii lor este salvat ca o singură valoare în index: $\text{norm}(d)$.

În continuare este prezentată formula practică de scor:

$$\text{score}(q, d) = \text{coord}(q, d) \times \text{queryNorm}(q) \times \sum_{t \in q} (tf(t, d) \times idf(t)^2 \times \text{boost}(t) \times \text{norm}(\text{field}(t), d)),$$

unde:

1. $tf(t, d)$ este corelat cu frecvența termenului în document. Documentele în care un termen apare de mai multe ori primesc un scor mai mare pentru acel termen. Lucene implementează astfel: $tf(t, d) = \sqrt{\text{freq}}$, unde freq reprezintă de câte ori apare termenul în document.
2. $idf(t)$ este inversul frecvenței termenului la nivel de index. Acest lucru înseamnă că termenii mai rari au o contribuție mai mare la scor. Implementarea Lucene este: $idf(t) = 1 + \log\left(\frac{\text{numDocs}}{\text{docFreq} + 1}\right)$, unde numDocs reprezintă numărul de documente din index și docFreq reprezintă numărul de documente în care apare termenul.
3. $\text{coord}(q, d)$ este o componentă calculată la momentul căutării: $\text{coord}(q, d) = \frac{\text{overlap}}{\text{maxOverlap}}$, unde overlap reprezintă numărul de termeni din interogare care se regăsesc în document și maxOverlap , numărul de termeni ai interogării.
4. $\text{queryNorm}(q)$ este factorul de normalizare folosit pentru a face scorurile pentru diferite interogări comparabile. Acest factor nu afectează ierarhizarea documentelor din moment ce este același pentru fiecare document. Implementarea implicită Lucene computează norma euclidiană a vectorului ponderilor (ajustate de nivelele de boost):

$$\text{queryNorm}(q) = \frac{1}{\sqrt{\text{boost}(q) \times \sum_{t \in q} (\text{idf}(t) \times \text{boost}(t))^2}}$$

5. $\text{boost}(t)$ reprezintă nivelul de boost al termenului; acesta poate fi setat din sintaxa interogării dacă este folosit parserul de interogări pus la dispoziție de Lucene, sau prin intermediul api-ului obiectului *Query*.

- 6.1.2 Pachetul Benchmark
- 6.1.3 Apache Open Relevance Project
- 6.2 Construirea framework-ului
- 6.3 Rezultate

Capitolul 7

Concluzii

Bibliografie

- [1] Christopher D. Manning, Prabhakar Raghavan and Hinrich Schutze, *Introduction to Information Retrieval*, Cambridge University Press. 2008.
- [2] http://en.wikipedia.org/wiki/Information_retrieval