

Universitatea "Babeş-Bolyai", Cluj-Napoca, România
Facultatea de Matematică și Informatică

Optimizarea Regăsirii Informației în Spațiul World Wide Web

Teză de doctorat

Student doctorand:
Ioan BĂDĂRÎNZĂ
Coordonator științific:
Prof. Dr. Florian Mircea BOIAN

2018

Rezumat

În această teză, cu titlul "Optimizarea Regăsirii Informației în Spațiul World Wide Web", am abordat un subiect de actualitate, care este intens studiat de către cercetători și anume regăsirea informației într-o colecție de date.

În zilele noastre, numărul de documente web din spațiul World Wide Web crește foarte repede. Această viteză de creștere a numărului de pagini web disponibile în Internet este datorată faptului că utilizatorii pot adăuga conținut rapid și ușor. Ei pot încărca, într-un mod foarte facil, multiple fișiere de date, de la fișiere text la cele audio și video și, de asemenea, pot edita și schimba conținutul paginilor în fiecare zi. Problema care apare în urma acestei creșteri rapide este legată de faptul că a devenit foarte greu pentru o persoană să găsească răspunsuri la anumite nevoi de informare în toate aceste date disponibile în Internet. Pentru a face ca aceste date să fie disponibile utilizatorilor, s-au făcut multe cercetări în domeniul de regăsire a informației. În 1958, Luhn menționează pentru prima dată problema de regăsire a informației, când a și propus prima modalitate de a determina utilitatea unui cuvânt pentru un anumit document din care acesta face parte. Această propunere a evoluat foarte mult în ultimele decenii. Regăsirea informației este un proces prin care un sistem software care primește ca date de intrare necesitatea de informare reală, exprimată printr-un query, returnează ca răspuns o listă de documente relevante, documente care au fost extrase dintr-o colecție foarte mare de date. Pentru a putea efectua o astfel de căutare într-un sistem de regăsire a informației, fiecărui cuvânt din document îi este atribuită o anumită pondere care determină cât de reprezentativ este acel cuvânt pentru acel document. Mai apoi, aceste ponderi sunt folosite în diverși algoritmi de clasare care returnează o lista de documente, ordonate după gradul de similaritate cu query-ul exprimat de utilizator.

Această teză se concentrează pe două subdomenii ale sistemelor de regăsire a informației: căutarea textuală bazată pe conceptul de "bag of words" și personalizarea sugestiilor de interogare. Conceptul de "bag of words" este cel mai utilizat concept în sistemele de regăsire a informației, dar are un dezavantaj major: sensul cuvintelor și implicit al propozițiilor și frazelor dintr-un document se pierde o dată cu indexarea lor. În acest sens, am făcut unele studii în care am implementat diferite instrumente, pe care le-am folosit cu scopul de a demonstra eficiența

propunerii noastre pentru a îmbunătății acest dezavantaj al sistemelor care implementează conceptul de "bag of words". O altă parte a sistemelor de regăsire a informației care ne-a captat atenția este mecanismul de sugestie a interogărilor. În ultimul timp, s-au făcut multe cercetări în acest domeniu, în special în personalizarea sugestiilor de interogare. Aceste sugestii sunt prezentate utilizatorului în care momentul în încercă să scrie query-ul prin care își exprimă necesitatea de informare. Am identificat faptul că niciunul din sistemele propuse de alți cercetători, nu folosește datele istorice personale de navigare a paginilor web, atunci când oferă lista de sugestii de interogare. În acest sens, am construit un instrument prin care am început să colectăm astfel de date despre utilizatori și am demonstrat că acest istoric personal de navigare al paginilor web poate fi folosit pentru a îmbunătății calitatea sugestiilor oferite utilizatorilor.

Cuvintele cheie relevante pentru subiectul abordat de această teză sunt: regăsirea informațiilor, căutare textuală, modele de clasificare, modelul boolean, modelul probabilistic, modelul spațiului vectorial, index, index inversat, bm25, index grupat, algoritm de clasare, interogare, sugestie interogare, sub-query, sugestie personalizată, analiză sintactică, index sintactic.

Teza are următoarea structură:

1 Introducere

1.1 Motivație

1.2 Enunțul Problemei

1.3 Tematica Tezei și Contribuțiile Principale

1.4 Structura Tezei

2 Fundamente și Concepte în Regăsirea Informatiei

2.1 Concepte. Fluxul de Lucru în Regăsirea Informatiei

2.1.1 Preprocesarea

2.1.1.1 Împărțirea în Termeni

2.1.1.2 Eliminarea Cuvintelor de Legătură

2.1.1.3 Reducerea la Forma de Bază

2.1.2 Compararea

2.2 Indexarea

2.2.1 Indexul de Incidentă Termen-Document

2.2.2 Indexul Inversat

2.2.3 Ponderea Termenilor

2.2.4 Crawling

-
- 2.3 Modele de Clasare
 - 2.3.1 Modelul Boolean
 - 2.3.2 Modelul Spațiului Vectorial
 - 2.3.3 Modelul Probabilistic
 - 2.3.3.1 Divergența de la Aleatoriu
 - 2.3.3.2 Modelul de Independență Binară și Principiul Probabilistic de Clasare
 - 2.3.3.3 BM25
 - 2.3.3.4 BM25F
 - 2.4 Sisteme de Sugestie
 - 2.4.1 Auto-Completarea Query-ului
 - 2.4.2 Sugestii de Interogare
 - 2.5 Evaluare
 - 2.5.1 Evaluarea Sistemelor de Regăsire a Informației
 - 2.5.2 Evaluarea Sugestiilor de Query-uri
- 3 Indecși Grupați pe Nivele în Regăsirea Informației Textuale
- 3.1 Indecși Grupați pe Nivele
 - 3.1.1 Preliminarii
 - 3.1.2 Structura Indexului Folosit în Sistem
 - 3.1.3 Algoritmul de Regăsire al Sistemului
 - 3.1.4 Evaluarea
 - 3.1.5 Concluzii
 - 3.2 Indecși Sintactici
 - 3.2.1 Indecși Sintactici în Regăsirea Informației și Procesarea Limbajului Natural
 - 3.2.2 Indecși Sintactici în Regăsirea Informației Textuale
 - 3.2.2.1 Metoda 1: Analiza Sintactică a Query-urilor de Tip Frază
 - 3.2.2.2 Metoda 2: Reducerea Mărimii Indexului Inversat
 - 3.2.2.3 Metoda 3: Indexul Sintactic
 - 3.2.3 Concluzii
 - 3.3 Ponderi Personalizate ale Termenilor
 - 3.3.1 Funcții de Clasare Folosite în Sistemele de Regăsire a Informației
 - 3.3.2 Ponderi Parțial Personalizate pentru Termenii din Query

3.3.3	Evaluarea
3.3.4	Concluzii
4	Personalizarea Sugestiilor de Interogare și Auto-Completărilor
4.1	Analiză a Istoriei de Navigare Web a Utilizatorului
4.1.1	Căutări Personalizate
4.1.2	Arhitectura Extensiei de Browser
4.1.3	Analiza Datelor Colectate
4.1.4	Concluzii
4.2	Personalizarea Sugestiilor de Interogare
4.2.1	Personalizarea Sugestiilor
4.2.2	Metoda Propusă pentru Personalizarea Sugestiilor de Interogare
4.2.3	Evaluarea
4.2.4	Concluzii
5	Concluzii și Direcții de Cercetare
5.1	Realizări
5.2	Direcții de Cercetare

În prima parte a tezei am descris datele generale despre regăsirea informației textuale. La început am prezentat o scurtă istorie a regăsirii de informații într-o colecție de date, când a fost pusă problema pentru prima dată și cum a ajuns să fie atât de bine cunoscută. Am continuat cu prezentarea atât a conceptelor fundamentale cât și a modului de lucru al unui sistem de regăsire a informației. Mai întâi am discutat despre componentele unui sistem de regăsire a informației și de ce este important pasul de preprocesare, după care am analizat cum este construit indexul unui sistem de regăsire a informației, ca mai apoi să descriem detaliile despre diferite modele de clasare, cum ar fi modelul spațiului vectorial, modelul probabilistic și modelul boolean. După acestea, am discutat despre o caracteristică mai nouă a sistemelor de regăsire a informației și anume sugerarea de query-uri în timp ce utilizatorul începe să tasteze query-ul dorit. Această caracteristică ajută utilizatorul să își definească mai bine query-ul prin care dorește să își exprime necesitatea de găsire a unor informații. În încheiere, am descris modalitățile de evaluare a sistemelor de regăsire a informației și a rezultatelor returnate de acesta și, nu în ultimul rând, cum putem evalua lista de sugestii returnată de un sistem de sugestii pentru un anumit query.

În continuare am prezentat contribuțiile majore din această teză:

Capitolul 3.1 - Indecși Grupați pe Nivele: În acest capitol este prezentat un sistem de regăsire a informației bazat pe indecși grupați pe nivele, care are la bază modelul spațiului vectorial. Acest sistem poate fi utilizat atât pentru căutarea exactă de fraze, cât și pentru o căutare mai bună a cuvintelor cheie prin utilizarea ponderii de proximitate în funcția de măsurare a similarității. Procesul de regăsire a informației este bazat pe grupuri de documente care sunt similare între ele, iar în momentul căutării sunt verificați doar liderii acestor grupuri, nu și toate documentele din baza de date.

Capitolul 3.2 - Indecși Sintactici: Acest capitol descrie trei metode prin care se pot încorpora metadate într-un sistem de regăsire a informației textuale. Prima metodă implică doar o analiză sintactică a query-ului, iar pe baza categoriei gramaticale generează o pondere diferită pentru fiecare termen din query, ponderi care mai apoi vor fi folosite în algoritmul de căutare. A doua metodă implică optimizarea spațiului de stocare a indexului inversat prin stocarea doar a acelor termeni care apar ca subiect și predicat într-un document. A treia metodă propune construirea unui index complet sintactic, ceea ce înseamnă că pentru fiecare termen din colecție, indexul inversat va stoca și categoria gramaticală a termenului, alături de frecvența termenului într-un document și frecvența inversă a termenului.

Capitolul 3.3 - Ponderi Personalizate ale Termenilor: După efectuarea unei analize a valorilor IDF pentru termenii din colecția Reuters, am observat că pentru interogările care conțin perechi de termeni, unul cu valori foarte mici ale IDF-ului și altul cu valori IDF foarte mari, sistemele de regăsire a informației nu returnează documente relevante. În acest capitol analizăm o metodă bazată pe modelul probailistic care permite utilizatorului să modifice ponderile termenilor din query în așa fel încât sistemele de regăsire a informației să returneze rezultate mai relevante. Această metodă nu are rolul de a spori relevanța rezultatelor pentru căutările generale, ci are scopul de a crește relevanța rezultatelor returnate pentru anumite căutări specifice în care termenii query-ului folosit au valori IDF disproporționate.

Capitolul 4.1 - Analiză a Istoriei de Navigare Web a Utilizatorului: În acest capitol subliniem rezultatele cercetării în sistemele de sugestii de interogare cu un accent deosebit acordat personalizării acestora. Sugestiile de interogare sunt un mecanism prin care utilizatorul primește o listă de interogări posibile de unde poate alege și poate efectua o căutare într-un mod mult mai rapid și mai ușor. Am evaluat utilitatea istoricului de navigare a paginilor web a unui utilizator în generarea de noi sugestii. Timp de o luna de zile am realizat un experiment în care, cu ajutorul unui plugin de Chrome, am colectat informații despre paginile web vizitate de utilizatori, alături de interogările pe care aceștia le făceau pe motorul de căutare Google și am constatat că aproximativ 32% din aceste interogări ar fi putut fi prezise din acea istorie de navigare a fiecărui utilizator în parte.

Capitolul 4.2 - Personalizarea Sugestiilor de Interogare: După ce am observat

faptul că istoricul personal de navigare al paginilor web poate fi folosit pentru a prezice următorul query pe care utilizatorul va încerca să îl folosească, în acest capitol am propus o nouă metodă de personalizare a ordinii în care sunt prezentate sugestiile de interogare. Astfel, sugestiile mai relevante pentru utilizator vor fi plasate pe un loc din listă cât mai sus (să aibă un rang mai mare). Modalitatea prin care am realizat acest lucru este definirea unui scor *Personal Temporal Query Suggestion* pentru fiecare sugestie din listă. Acest scor ține cont de istoricul personal, foarte recent și scurt, de navigare a paginilor web, care este folosit în continuare pentru a calcula un *HybridPageScore*. Capitolul mai conține evaluarea și experimentele pe care le-am efectuat pentru metoda propusă.

În capitolul de concluzii, am prezentat principalele contribuții ale acestei teze la optimizarea sistemelor de regăsire a informațiilor într-o colecție de date și direcțiile viitoare de cercetare în acest domeniu.