

Sisteme de regăsire a informațiilor – Metode de indexare și căutare a datelor

Anca Farcaș

Rezumat

Tema abordată în lucrarea de diplomă este **procesul de indexare și căutare în motoarele de căutare**. Indexarea poate fi privită ca un proces ce oferă o relație de ordine asupra datelor pe care le analizăm. În această lucrare se va pune accent pe două tipuri de indexare și anume, **indexarea directă și indexare inversă**.

Indexarea reprezintă procesul prin care datele unei colecții sunt reorganizate în scopul de a fi regăsite ușor și precis. **Indexarea directă** are ca scop determinarea indecșilor relativ la document. **Indexare inversă** are drept scop determinarea documentelor relativ la index. Aceasta reprezintă componenta de bază a oricărui motor de căutare, utilă în determinarea relevanței unui cuvânt la colecția de documente în care se face căutarea.

Printre etapele principale ale indexării directe, se numără **pre-procesarea** și construirea **indecșilor adecvați**. Pentru preprocesare, a fost considerată componenta de analiza a formei de bază a unui cuvânt. Pentru a aduce cuvântul la forma sa canonică au fost construiți doi arbori lingvistici care comunică între ei, unul din aceștia fiind populați cu toate cuvintele limbii engleze, iar al doilea cu formele neregulate a substantivelor, verbelor, adverbilor (exemplu: man - men, mouse - mice), fiecare cuvânt din arborele neregulat având legătură cu forma de bază a cuvântului ce îi coincide în arborele dicționar. Scopul acestei abordări este de a soluționa acele cazuri în care algoritmi clasici de stemming (*Porter, Lovins și Paice*) nu oferă rezultate satisfăcătoare. Această formă de preprocesare este aplicată atât în faza de construire efectivă a indexului, cât și în componenta de căutare.

Fiecare fișier din colecția de documente deja indexate direct, vor fi trecute prin procesul de **indexare inversă**, proces care va ajuta la cautarea rapidă a unui cuvânt în întreaga colecție de fișiere.

Pentru componenta de căutare, accentul s-a pus pe **modul de căutare** de tip **boolean** și cel bazat pe **reprezentarea vectorială**.

În ceea ce privește **modelul boolean** de căutare, reprezentarea documentului se realizează sub formă de ponderi. Termenii interogării sunt combinați logic utilizând operatorii booleani AND, OR și NOT, iar regăsirea documentelor căutate se bazează aritmetica mulțimilor și a criteriului deciziei binare. Acest modul de căutare boolean a ajutat la restrângerea colecției de documente pentru interogarea utilizatorului.

Pentru **modulul de căutare bazată pe reprezentarea vectorială**, reprezentarea documentului, cât și a interogării, constă în vectori de ponderi și se obține rapid dintr-un **index direct cantitativ**. Regăsirea documentelor de interes se bazează pe distanța dintre două documente calculată ca fiind unghiul dintre doi vectori n-dimensionali. Acest modulul, bazat pe reprezentarea vectorială a oferit seturi relevante de documente pentru cererea utilizatorului.