

UNIVERSITATEA TEHNICĂ „Gheorghe Asachi” din IAȘI
FACULTATEA DE AUTOMATICĂ ȘI CALCULATOARE
DOMENIUL: Calculatoare și tehnologia informației
SPECIALIZAREA: Tehnologia informației

Proiect RIW

Student
Andrișan Ionuț-Cosmin

Iași, 2020

1. Cerința proiectului

Implementați o aplicație care să realizeze funcția de căutare a unui sistem de regăsire a informațiilor dintr-un set de documente (txt, html).

Scenariul de test pentru proiect va fi următorul:

- aplicația va primi ca mesaj de intrare numele unui director ce conține un set de fișiere tip txt;
- acest director va trebui parcurs, recursiv, pentru a identifica acele fișiere txt și pentru a determina cele două forme de indexare indicate;
- modul de căutare ar trebui expus printr-o pagină HTML simplă (dar nu este obligatoriu – puteți opta pentru a dezvolta orice altă formă de interfață, inclusiv una de tip command line);

Încercați să realizați o implementare paralelă/distribuită pentru componentele ce implementează modulele de intrare și/sau căutare. Ca principal model de lucru, puteți porni de la modelul Map & Reduce.

Încercați să identificați și să analizați și alți algoritmi de stemming.

2. Etapele căutării

2.1. Construirea index-ului direct cantitativ

Aplicația primește ca intrare un director pe care îl parcurgem pentru a identifica fișierele pentru prelucrare. Directorul este parcurs astfel încât să fie identificate și fișierele din subdirectoare (în cazul în care în director avem subdirectoare).

Apoi, lista de fișiere obținută se sparge în cuvinte (pe care le stocăm într-un HashMap), contorizând în același timp și numărul de apariții a fiecărui cuvânt, astfel:

- 1) cuvântul din text, care a fost determinat în cadrul iterației curente, va fi testat contra unei liste de excepții – un dicționar al unei limbi nu conține, de exemplu, nume proprii; dacă acest cuvânt se regăsește în lista de excepții, atunci se va trece la următoarele etape de procesare (contorizare număr de apariții, etc.);
- 2) cuvântul din text, care a fost determinat în cadrul iterației curente, va fi testat contra unei liste de stop-word-uri – cuvintele de legătură care în mod uzual nu aduc informații noi pentru motoarele de căutare; dacă acest cuvânt curent determinat se regăsește într-o astfel de listă de stop-word-uri, atunci acesta va fi eliminat din procesările ulterioare;
- 3) cuvântului din text, care a fost determinat în cadrul iterației curente, i se va aplica un algoritm de stemming (algoritmul lui Porter) pentru a se ajunge la o formă de bază a cuvântului, numită și formă canonică.

Stocarea indexului direct cantitativ se face astfel: se stochează cheia (documentele indexate), cât și valorile asociate (cuvintele din cadrul documentelor și numărul de apariții a fiecăruia) cheilor. Acest index se salvează atât într-un fișier txt, cât și într-o bază de date non-relațională (MongoDB).

2.2. Construirea index-ului indirect cantitativ

Index-ul indirect cantitativ se obține pe baza index-ului direct cantitativ și este de forma: cheia (cuvântul indexat) și documentele în care găsim aceasta cheie (aceasta reprezintă fișierele în care se găsește cuvântul, cât și numărul de apariții a aceluia cuvânt în documentul respectiv). De asemenea, acest index se salvează atât într-un fișier txt, cât și într-o bază de date non-relațională (MongoDB).

3. Modelul de căutare booleană

Reprezentarea interogării – termenii interogării (sau cheile de căutare) sunt combinate logic utilizând operatorii booleani AND, OR și/sau NOT.

Regăsirea documentului – se bazează criteriul deciziei binare și pe aritmetica mulțimilor.

Avantaje:

- Este un model de cautare simplu, cu un formalism bine pus la punct, neambiguu.
- Poate fi implementat ușor și poate raspunde rapid pentru interogările uzuale ale utilizatorilor.

De zavantaje:

- Datorită simplității, este un model foarte rigid.
- Interogările complexe nu pot fi realizate direct.
- Nu poate fi controlată cu exactitate dimensiunea exactă a răspunsului.
- Nu ofera un mecanism direct de feedback din partea utilizatorilor.

Principalii pași implicați:

1. se citește interogarea utilizator;
2. se izolează operanzii (cuvintele) de operatori;
3. cuvintele se procesează conform modelului utilizat în construirea index-ului corespunzător;
4. se izolează pe baza index-ului invers, pentru fiecare cuvânt în parte lista de documente ce conțin termenul respectiv;
5. se realizează, rând pe rând, operațiile indicate **AND**, **OR** și/sau **NOT**:
 - **AND** echivalează cu intersecția a două mulțimi;
 - **OR** echivalează cu reuniunea a două mulțimi;
 - **NOT** echivalează cu diferența dintre două mulțimi;
6. rezultatul obținut este prezentat utilizatorului.

4. Procesarea cuvintelor – algoritmul lui Porter (limba engleza)

Stemmingul lui Porter este un algoritm care este folosit pentru aducerea cuvintelor din limba engleză într-o formă canonică. Algoritmul este văzut ca un algoritm „înghețat”, care este definit strict, nu poate fi modificat și produce un număr mare de erori.

De multe ori, algoritmul lui Porter pentru stemming nu aduce cuvântul la forma canonică, ci la o formă greșită. Totuși, motivul stemmer-ului este de a aduce o serie de cuvinte la aceeași formă.