

Laborator nr. 4

Modul de căutare booleană

an univ. 2018 – 2019

1 Scurt breviar teoretic

1.1 Caracteristici

Reprezentarea documentului – vector de ponderi, pentru care ponderea w_{ij} este definită conform ecuației de mai jos:

$$w_{ij} = \begin{cases} 1 & \text{dacă } t_i \in d_j \\ 0 & \text{altfel} \end{cases} \quad (1)$$

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

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

Avantaje:

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

Dezavantaje

- 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 oferă un mecanism direct de *feedback* din partea utilizatorilor.

1.2 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.

2 Aplicație propusă

Implementați o aplicație care să realizeze funcția de căutare a unui sistem de regăsire de informații conform **modelului boolean** ([1], *Capitolul 1: Boolean retrieval*). Operațiile dorite sunt:

- $k1 \text{ AND } k2$: setul de documente care conțin atât termenul $k1$, cât și termenul $k2$ – **intersecție de mulțimi**;
- $k1 \text{ OR } k2$: setul de documente care conțin fie termenul $k1$, fie termenul $k2$ – **reuniune de mulțimi**;
- $k1 \text{ NOT } k2$: setul de documente care conțin termenul $k1$, dar nu conțin termenul $k2$ – **diferență de mulțimi**.

Aplicația trebuie să poată procesa și interogări ce conțin mai mult de două chei de căutare (exemplu: $k1 \text{ AND } k2 \text{ NOT } k3 \text{ OR } k4$). Indexul invers de test este cel determinat în cadrul laboratorului nr 3.

Bibliografie

- [1] Christopher D. Manning et. al. Introduction to Information Retrieval. <http://nlp.stanford.edu/IR-book/information-retrieval-book.html>, 2009.