

Regăsirea Informațiilor pe Web

P01 – Indexare și căutare

Student: Neștian Mihai

1. Descriere

Proiectul conține două părți: partea de indexare și cea de căutare.

Modulul de indexare primește ca intrare calea și numele unui director ce conține un set de fișiere de tip *.txt*. Acest director este parcurs recursiv și pentru fiecare fișier este construit *indexul direct*, stocat apoi într-o bază de date *MongoDB*. La construirea indexului direct, este parcurs fiecare fișier caracter după caracter și sunt determinate cuvintele. Fiecare cuvânt trece prin două teste înainte de a fi adăugat în index. În primul rând se verifică dacă face parte din lista de excepții (din fișierul *exceptions*). Dacă nu este o excepție, se verifică dacă se găsește în lista de cuvinte ignorate (din fișierul *stopwords*). Dacă trece aceste teste, cuvântul este adus la forma canonică prin trecerea lui prin algoritmul *Porter Stemming* din biblioteca *nlk* și este adăugat apoi în index, împreună cu numărul de apariții în document. Pe baza indexului direct, se construiește *indexul indirect cantitativ*, care este stocat, la fel ca cel direct, în baza de date *MongoDB*.

Am ales să folosesc algoritmul lui Porter deoarece are o acuratețe mai mare față de cel al lui Lancaster. Timpul de prelucrare este același, prelucrarea fișierelor a durat 12 secunde cu ambii algoritmi. Am mai încercat și un lemmatizer. Acuratețea acestuia este mult mai bună față de cea a algoritmilor de stemming dar timpul de prelucrare este mult mai mare, prelucrarea celor 26 de fișiere a durat 1:53 față de 12 secunde pentru algoritmi de stemming.

Cuvant	Porter	Lancaster	Lemmatizer
programs	program	program	program
programming	program	program	program
programmers	programm	program	programmer
languages	languag	langu	language
peoples	peopl	peopl	people
programmer	programm	program	programmer
involving	involv	involv	involve
swimming	swim	swim	swim
spliting	splite	spliting	split
are	are	ar	be
going	go	going	go
tissues	tissu	tissu	tissue
bigger	bigger	big	big
smallest	smallest	smallest	small
churches	church	church	church
biggest	biggest	biggest	big
Churchill	churchil	churchil	Churchill

Modulul de căutare cuprinde două metode: căutarea booleană și căutarea vectorială. Acesta primește ca intrare o expresie de căutare din partea utilizatorului și returnează numele fișierelor relevante pentru expresia introdusă și scorul obținut de acestea. Pentru căutarea booleană această expresie trebuie să conțină operatori booleani între cuvinte (*AND*, *OR*, *NOT*) care sunt reprezentați de (&, |, !). Căutarea se face pe baza indecsului indirect, iar criteriile de relevanță sunt: similaritatea cosinus pentru căutarea vectorială și criteriul deciziei binare și logica mulțimilor pentru căutarea booleană.

Pentru stocarea datelor s-a folosit o bază de date MongoDB, în care sunt stocați indecșii direct și indirect, dar și matricea de ponderi folosită la căutarea booleană, *inverse document frequency* și valorile ponderilor documentelor (*term frequency * inverse document frequency*) folosite pentru căutarea vectorială.

Toate fișierele de test se găsesc în directorul *Structura*. Cele generate de program se găsesc în *Direct_Index* și *Indirect_Index*.

2. Utilizare

Aplicația folosește o bază de date MongoDB, configurată cu setările implicite (*localhost*, port 27017). Numele bazei de date folosite este *mydatabase* și conține un număr variabil de colecții, în funcție de numărul fișierelor de intrare:

- ID 1-N - aici este stocat indexul direct
- *direct_index* – aici este stocată maparea fișierelor, ce colecție de index îi revine fiecărui fișier
- *indirect_index_cantitativ* – aici este stocat indexul invers
- *info_collection* – cuprinde TF, IDF și vectorul asociat fiecărui document

Aplicația conține și o interfață grafică. Interfața este de tip *tkinter* și conține un câsuță text în care utilizatorul poate introduce expresia de căutare, un buton: *Search și încă o căsuța text unde vor fi afișate documentele rezultate în urma căutării vectoriale*. Prin apăsarea butonului Search se va apela funcția de căutare vectorială.

3. Auto-evaluare

baza - **1 punct**

1. parcurgere recursivă structura de directoare – **0.5 puncte**

2. procesare cuvinte + stemming – **2 puncte**

3. indexare indirectă - **3 puncte**

4. căutare (în indecși preîncărcați)

vectorială (cu restrângerea colecției de lucru) - **3 puncte**

5. stocare indecși în mongodb - **1 punct**

Total: 9.5 puncte + 1 punct bonus