

Proiect 1

Indexare și căutare

an univ. 2018 – 2019

1 Prezentarea temei de proiect

Prima componentă de proiect a disciplinei „Regăsirea Informațiilor pe WEB” vă propune să realizați un set de modificări asupra aplicațiilor dezvoltate pe parcursul laboratoarelor 1 – 4 astfel încât:

1. să obțineți o procesare adecvată a cuvintelor determinate în cadrul unui text;
2. să studiați și alte tipuri de baze de date și mecanisme de stocare de date (precum MongoDB);
3. să studiați alte metode de realizare a operațiilor de căutare.

Astfel, pentru prima cerință de mai sus, metodele de construire a indecșilor direcți și inverși trebuie să implementeze mecanisme prin intermediul cărora un cuvânt de dicționar să fie adus la așa numita *formă canonică* [1] (capitolul 2). Indecșii astfel determinați, vor fi apoi stocați prin intermediul unor baze de date ne-relaționale, precum MongoDB.

După cum a fost amintit în cadrul laboratorului nr. 4, căutarea booleană nu include mecanisme de determinare a unui eventual scor de relevanță pentru rezultatele identificate. A treia cerință vă propune să studiați impactul *distanței de tip cosinus* asupra rezultatelor obținute de un motor de căutare.

2 Scurt breviar teoretic

2.1 Procesarea cuvintelor – algoritmul lui Porter (limba engleză)

Procesele prin care un cuvânt este adus la forma sa de bază poartă numele de *stemming* și *lematizare* [1]. Cele două modele de procesare urmăresc același scop: reducerea dimensiunii indecșilor. Pentru cuvintele ce aparțin de dicționarul unei limbi nu are sens stocarea tuturor formelor sub care aceste cuvinte pot fi derivate/conjugate (deoarece aceste transformări nu aduc efectiv informații noi procesului de indexare).

Modelele de analiză bazate pe *lematizare* sunt axate pe analiză morfologică și, cel puțin teoretic, ar trebui să ofere rezultate cu un grad ridicat de precizie. Totuși, procesele sunt lente și considerabil mai greoi de implementat.

Procesele de tip *stemming* se bazează pe trunchierea cuvintelor (de exemplu, în limba engleză pluralul cuvintelor se obține în majoritatea cazurilor, prin adăugarea caracterului 's'). Aceste tehnici sunt considerabil mai simple, cu un timp de răspuns mult redus, dar nu oferă aceeași precizie. Unul dintre cei mai renumiți algoritmi este algoritmul lui Porter – pentru detalii vezi [2].

2.2 Căutare vectorială

Etapele de lucru implicate de această formă de căutare sunt descrise de algoritmul 1.

Algoritm 1 cautareVectoriala(D, q)

- 1: transforma fiecare document $d_j \in D$ în $\vec{d}_j = \{key : tf(key, d) \cdot idf(key)\}$
 - 2: transforma interogarea $\vec{q} = \{key : tf(key, d) \cdot idf(key)\}$
 - 3: **for all** $\vec{d}_j \in D$ **do**
 - 4: calculeaza $s_j = similaritateCosinus(\vec{d}_j, \vec{q})$
 - 5: **end for**
 - 6: sorteaza documentele descrescator din punct de vedere al scorului anterior s_j
 - 7: **return** setul relevant de documente
-

Observație: Coeficienții tf și idf sunt descriși în cursul 3, slide 15, formulele (2) și (3).

2.3 Tutorial MongoDB

NoSQL (baze de date nerelaționale) este un model care permite memorarea unor volume imense de date, fără o structură fixă și care se poate schimba, care oferă scalare pe orizontală [7]. MongoDB este o baza de date de documente, open-source, high performance, high availability, ce ofera scalare automata [3]. Documentul în MongoDB este o structură de date formată din perechi cheie valoare, similară cu obiectele JSON.

Exemplu stocare index direct și invers într-o bază de date mongodb
<p>Doc1: Data mining este o tehnica noua de analiza a datelor. Doc2: Tehnicile data mining pot aduce informatii noi. Doc3: Datele sunt colectate prin tehnici specifice.</p>
<p><i>DB indexdirectcantitativ</i></p> <pre>{ "doc": "Doc1", "terms": [{ "t": "analiza", "c": 1 }, { "t": "data", "c": 2 }, { "t": "mining", "c": 1 }, { "t": "noutate", "c": 1 }, { "t": "tehnica", "c": 1 }] }</pre> <pre>{ "doc": "Doc2", "terms": [{ "t": "data", "c": 1 }, { "t": "informatie", "c": 2 }, { "t": "mining", "c": 1 }, { "t": "noutate", "c": 1 }, { "t": "tehnica", "c": 1 }] }</pre> <pre>{ "doc": "Doc3", "terms": [{ "t": "data", "c": 1 }, { "t": "tehnica", "c": 1 }] }</pre>
<p><i>DB indexinverscantitativ</i></p> <pre>{ "term": "analiza", "docs": [{ "d": "Doc1", "c": 1 }] }</pre> <pre>{ "term": "data", "docs": [{ "d": "Doc1", "c": 1 }, { "d": "Doc2", "c": 1 }, { "d": "Doc3", "c": 1 }] }</pre> <pre>{ "term": "informatie", "docs": [{ "d": "Doc2", "c": 1 }] }</pre> <pre>{ "term": "mining", "docs": [{ "d": "Doc1", "c": 1 }, { "d": "Doc2", "c": 1 }] }</pre> <pre>{ "term": "noutate", "docs": [{ "d": "Doc1", "c": 1 }, { "d": "Doc2", "c": 1 }] }</pre> <pre>{ "term": "tehnica", "docs": [{ "d": "Doc1", "c": 1 }, { "d": "Doc2", "c": 1 }, { "d": "Doc3", "c": 1 }] }</pre>
<p><i>Exemple de căutări:</i></p> <pre>db.indexinverscantitativ.find({ "term": "mining" })</pre> <pre>db.indexinverscantitativ.find({ \$or: [{ "term": "mining" }, { "term": "tehnica" }] })</pre>

Operații în mongo (inserări, căutări, actualizări): exemple pot fi consultate la adresa disponibilă în [4].
Java driver docs: [5].

```
1 MongoClient mongoClient = new MongoClient( "localhost" , 27017 );
2 MongoDB database = mongoClient.getDatabase("mydb");
3 MongoClient<Document> collection = database.getCollection("test");
4 Document doc = new Document("name", "MongoDB")
5     .append("type", "database")
6     .append("count", 1)
7     .append("info", new Document("x", 203).append("y", 102));
8 collection.insertOne(doc);
9
10 Document myDoc = collection.find(eq("i", 71)).first();
11 System.out.println(myDoc.toJson());
```

Pachetul Java corespunzător: [click aici](#).
Limbajul Python: [6].

3 Observații finale

- Scenariul de test pentru proiectele voastre va fi următorul:
 - aplicațiile vor primi ca mesaj de intrare numele unui director ce conține un set de fișiere de tip *txt*;
 - acest director va trebui parcurs, recursiv, pentru a identifica acele fișiere *txt* și pentru a determina cele două forme de indexare indicate;
 - modulul 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 implemeneză modulele de indexare și/sau cele de căutare. Ca principal model de lucru, puteți porni de la modelul **Map & Reduce**.
- Încercați să identificați și să analizați și alți algoritmi de *stemming*.

Bibliografie

- [1] Christopher D. Manning et. al. Introduction to Information Retrieval. <http://nlp.stanford.edu/IR-book/information-retrieval-book.html>, 2009.
- [2] Martin Porter. The Porter stemming algorithm. <http://snowball.tartarus.org/algorithms/porter/stemmer.html>.
- [3] mongoDB. Introduction to MongoDB. <https://docs.mongodb.com/manual/introduction/>.
- [4] mongoDB. MongoDB CRUD Operations. <https://docs.mongodb.com/manual/crud/>.
- [5] mongoDB. MongoDB Driver Quick Tour. <http://mongodb.github.io/mongo-java-driver/3.0/driver/getting-started/quick-tour/>.
- [6] mongoDB. Python MongoDB Drivers. <https://docs.mongodb.com/ecosystem/drivers/python/>.
- [7] mongoDB. What is NoSQL. <https://www.mongodb.com/nosql-explained>.

Barem evaluare proiect

Proiectele pornesc ca bază de notare de la **1 punct**.

Criteriu	Punctaj
1. Parcurgerea directorului cu setul de fișiere de indexat - <i>obligatoriu</i> parcurgerea unui singur nivel <i>sau</i> parcurgerea iterativă a unei structuri de directoare (inclusiv sub-directoarele)	0.25 puncte 1 punct
2. Preprocesare aplicarea (și înțelegerea modului de funcționare) unui algoritm de stemming pentru determinarea formelor canonice ale cuvintelor determinate	1 punct
3. Indexarea indirectă - <i>obligatoriu</i> (inclusiv generarea corectă a fișierlor de indecși)	4 puncte
4. Căutarea în indecșii determinați preîncărcați - <i>obligatoriu</i> doar căutare booleană <i>sau</i> căutare vectorială (cu restrângerea colecției de lucru)	1.5 puncte 3 puncte

Observații

Studentii vor fi rugați să completeze o auto-evaluare a proiectului dezvoltat. Auto-evaluările care nu diferă de nota acordată de titularii de laborator vor primi **suplimentar 0.5 puncte**. În plus, implementările pot atrage punctaj suplimentar (puncte bonus) pentru media finală a disciplinei, astfel:

stocarea indecșilor în colecții de documente de tip MongoDB – 2 puncte;

paralelizări/distribuiți eficiente ale volumului de lucru (eventual urmărind modelul MapReduce) – până la 3 puncte.