

Sistem de Regăsirea Informației pentru Limba Romana

Information Retrieval & Text Mining

Student: Lupașcu Marian

CLASELE

Proiectul conține următoarele clase:

- Indexer – clasa care pornește de la un set de documente pe care îl preia ca parametru în funcția main (args[0]) și creează un “inverted index” pe care îl salvează în folderul “.\index”. Aceasta clasa citește documentele cu ajutorul clasei DocumentReader (care folosește Tika) apoi se salvează ca documente txt (care conțin fix aceeași informație) într-un folder temporar pe baza căruia se construiește “inverted index-ul”, apoi folderul temporar este șters și informația este salvată pe disc în “inverted index”.
- Searcher – clasa care pornește de la “inverted index-ul” creat anterior și de la un string de căutare numit în continuare query. Aceasta clasa returnează documentele care sunt relevante pentru stringul de căutare pe baza unui scor de confidență.
- DocumentReader – clasa care parează diferite tipuri de documente cu Tika (txt, doc/docx sau pdf) și returnează informația ca plain-text.
- MyRomanianAnalyzer – clasa de procesare a informațiilor salvate de Indexer și a query-ului de căutare. Aceasta clasa este o extensie a clasei RomanianAnalyzer din Lucene (nu a putut fi moștenită din RomanianAnalyzer întrucât aceasta este clasa finală în Lucene), peste care au fost adăugate câteva features: metoda modifyStopWords(care îmbogățește lista default de stopwords din Lucene cu stopwords fără diacritice – dacă există stopword-ul câteva acum există și câteva) și în funcția de bază createComponents se mai adaugă câteva procesări. Acestea sunt eliminarea de diacritice a token-ilor apoi eliminarea de stopwords cu și fără diacritice, de unde rezultă doar tokeni fără stopwords, apoi ca stemming se folosește SnowballFilter și alți câțiva filtri de normalizare din Lucene.
- RemoveDiacriticalsFilter – o clasa filtru de eliminare de diacritice din tokeni, necesară pentru o preprocesare în MyRomanianAnalyzer.

RULARE

În continuare toate operațiile vor fi făcute din folderul de bază și anume P1.

1. Se creează un folder în P1 în care se adaugă documentele dorite să fie indexate. By default există folderul docs care conține 13 documente doc/x, 6 documente pdf, 4 documente txt și 2 imagini (deci în total 23 de documente ce vor fi analizate din 25 – cele 2 imagini nu vor fi analizate).
2. Se creează folderul dependencies în care se adaugă dependențele proiectului:
 - lucene-analyzers-common-8.6.3.jar,
 - lucene-core-8.6.3.jar
 - lucene-queryparser-8.6.3.jar
 - pdfbox-app-2.0.21.jar
 - tika-app-1.24.1.jar
3. Se pornește un terminal în P1 apoi se adaugă comanda de indexare a documentelor

```
java -Dfile.encoding=UTF-8 -classpath ".\out\production\P1;\dependencies\lucene-core-8.6.3.jar;\dependencies\tika-app-1.24.1.jar;\dependencies\pdfbox-app-2.0.21.jar;\dependencies\lucene-queryparser-8.6.3.jar;\dependencies\lucene-analyzers-common-8.6.3.jar" com.main.Indexer ".\docs"
```

4. După indexare se poate căuta diverse informații cu comanda. Proiectul a fost rulat cu Java 15.

```
java -Dfile.encoding=UTF-8 -classpath ".\out\production\P1;\dependencies\lucene-core-8.6.3.jar;\dependencies\tika-app-1.24.1.jar;\dependencies\pdfbox-app-2.0.21.jar;\dependencies\lucene-queryparser-8.6.3.jar;\dependencies\lucene-analyzers-common-8.6.3.jar" com.main.Searcher "de modificat"
```