

Recuperación de Información

Práctica 4. Indexación

Santiago Álvarez Valdivia

Ioannis Efthymiou

November 10th, 2023



Program Usage

This program (Practica4.jar) is meant to read information about Simpsons episodes (in a particular format) and create a Lucene index from the information.

An invocation of the program has three parameters: the operation to perform, the kind of data to index, and the directory where the csv files to be indexed are.

There are two operations available, **index** (which creates a new Lucene index from scratch) and **append** (which adds documents to an index created previously).

There are two data types known, **episodes** and **scripts**, which contain a summary of episode information and the complete information of what the characters of a Simpsons episode say, respectively.

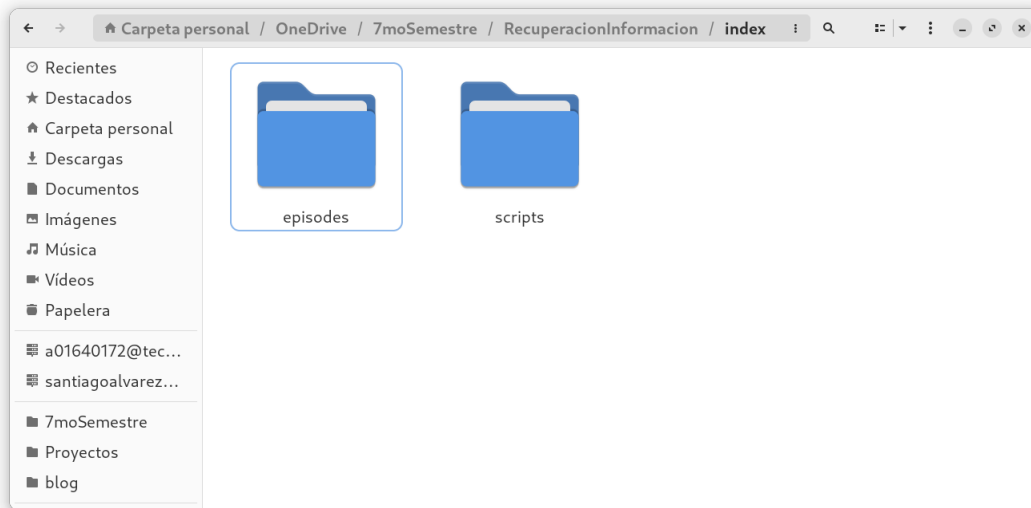
Then a typical invocation of the program would look like this:

```
java -jar Practica4.jar index episodes ./TestData/CapitulosUnidos
```

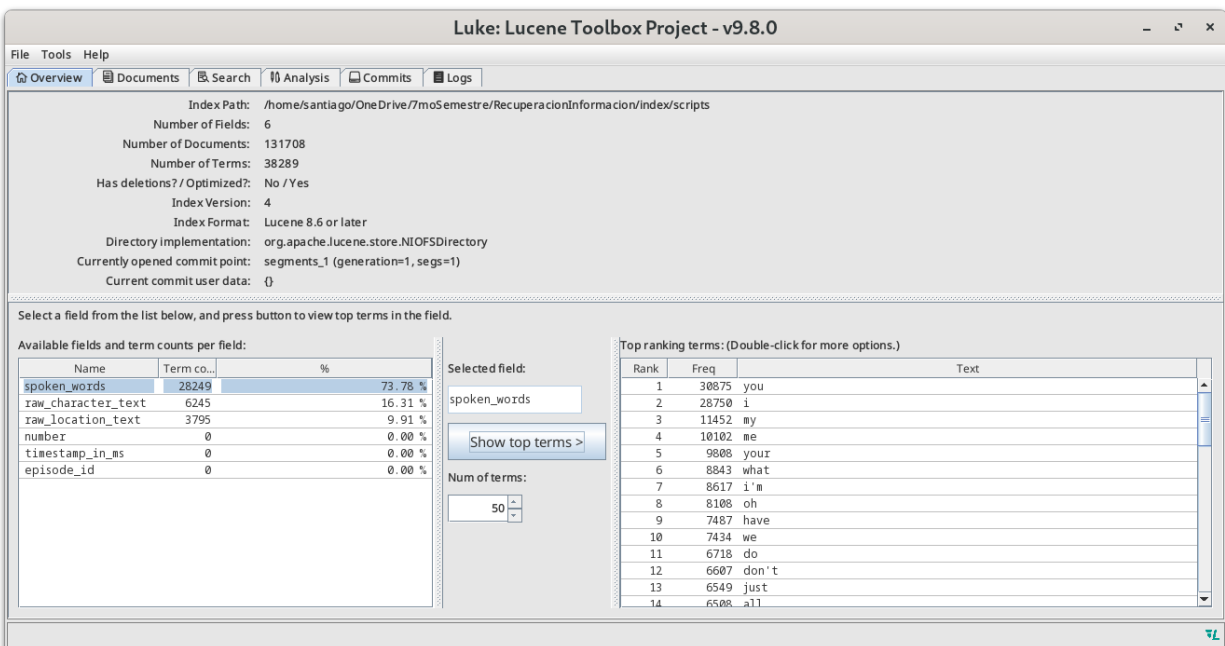
which would create a new index for the files in CapitulosUnidos, which are expected to have the format for episode summaries. In the case of the append operation, an index must already exist in the location where a new one would be created, otherwise we will get an error when the index is tried to be accessed.

Demo

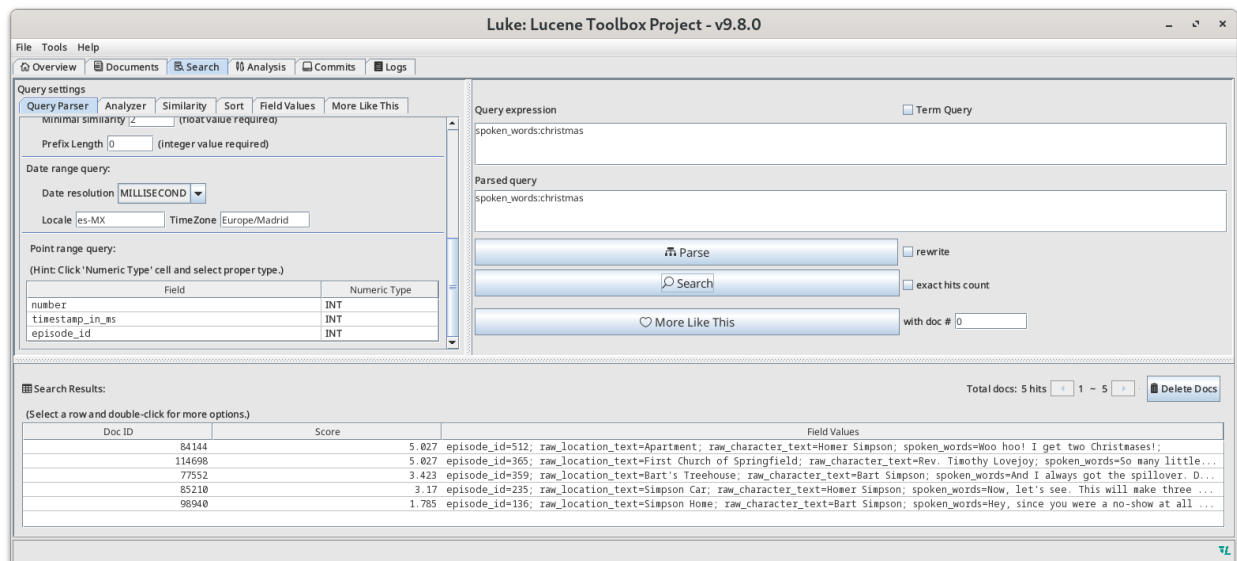
After running the program for both data types, we end up with this indexes



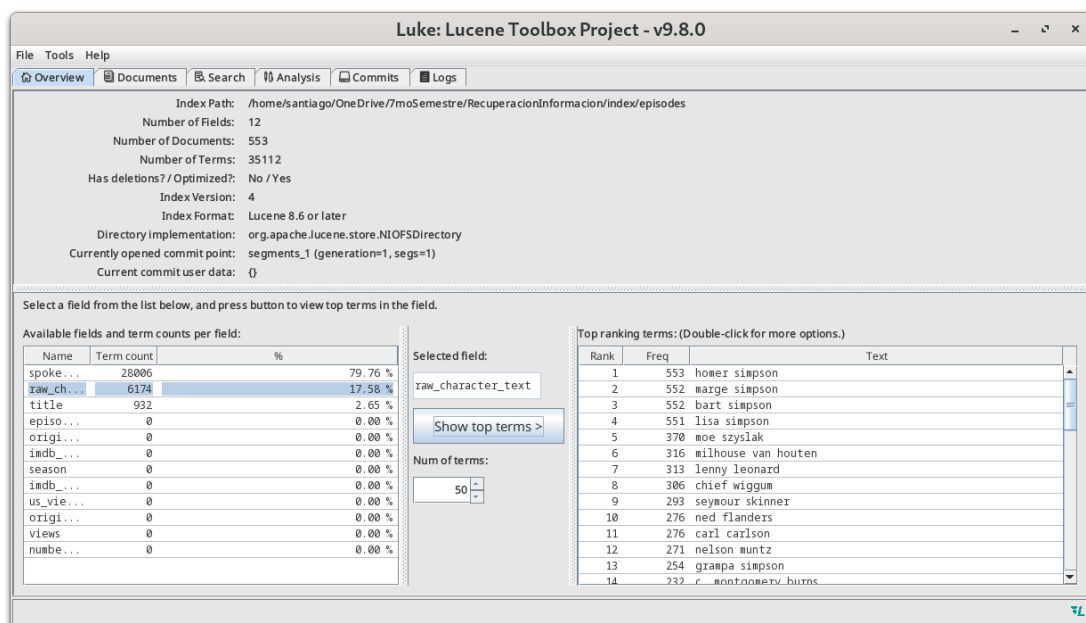
Opening the scripts index with Luke shows the indexed information



And we can perform searches on the indexed information



Meanwhile, this is how the episodes index looks like



We can also filter information on this index

Luke: Lucene Toolbox Project - v9.8.0

File Tools Help

Overview Documents Search Analysis Commits Logs

Query settings

Query Parser Analyzer Similarity Sort Field Values More Like This

Locale: es-MX TimeZone: Europe/Madrid

Point range query:
(Hint: Click 'Numeric Type' cell and select proper type.)

Field	Numeric Type
episode_id	INT
original_air_date	INT
imdb_votes	INT
season	INT
imdb_rating	INT
us_viewers_in_millions	INT
original_air_year	INT
views	INT
number_in_season	INT

Query expression

Season=3

Parsed query

Season:[3 TO 3]

Parse ☐ rewrite

Search ☐ exact hits count

More Like This with doc # 0

Search Results: Total docs: 24 hits 1 ~ 10 Delete Docs

(Select a row and double-click for more options.)

Doc ID	Score	Field Values
54	1	original_air_date=1992-02-06; title=Homer Alone; episode_id=50; us_viewers_in_millions=23.7; number_in_season=15; season=3; spoken_wor...
72	1	original_air_date=1992-03-12; title=Dog of Death; episode_id=54; us_viewers_in_millions=23.4; number_in_season=19; season=3; spoken_wor...
125	1	original_air_date=1991-10-17; title=Homer Defined; episode_id=48; us_viewers_in_millions=20.6; number_in_season=5; season=3; spoken_wor...
151	1	original_air_date=1992-05-07; title=Bart's Friend Falls in Love; episode_id=58; us_viewers_in_millions=19.5; number_in_season=23; season...
153	1	original_air_date=1992-02-20; title=Homer at the Bat; episode_id=52; us_viewers_in_millions=24.6; number_in_season=17; season=3; spoken...
167	1	original_air_date=1991-12-26; title=I Married Marjoe; episode_id=47; us_viewers_in_millions=21.9; number_in_season=12; season=3; spoken...
196	1	original_air_date=1992-02-27; title=Separate Vacations; episode_id=53; us_viewers_in_millions=23.7; number_in_season=18; season=3; spok...
231	1	original_air_date=1992-01-23; title=Lisa the Greek; episode_id=49; us_viewers_in_millions=23.2; number_in_season=14; season=3; spoken_w...
234	1	original_air_date=1992-02-13; title=Bart the Lover; episode_id=51; us_viewers_in_millions=20.5; number_in_season=16; season=3; spoken_w...
238	1	original_air_date=1991-11-14; title=Saturdays of Thunder; episode_id=44; us_viewers_in_millions=24.7; number_in_season=9; season=3; spo...