

UNIVERSIDAD DE GRANADA
E.T.S.I INFORMÁTICA Y TELECOMUNICACIÓN



Recuperación de Información (RI)

**Implementación de un Sistema de
Recuperación de Información utilizando
Lucene**

Parte A. Indexación

Curso 2017-2018

Cuarto Curso del Grado en Ingeniería Informática

María Camarero Granados

Javier Gómez Luzón

Francisco Porcel Molina

Índice

- Análisis previo de los requisitos.....3
 - Interfaz de búsqueda.....3
 - Analizador.....4
 - Facetas.....4
- Diseño de la solución.....5
- Manual de usuario.....7

Análisis previo de los requisitos

Se nos pide implementar un programa que permita añadir documentos a un índice creado utilizando la biblioteca Lucene, que posteriormente formara parte de un Sistema de Recuperación de Información que servirá para buscar artículos científicos.

Además, se debe utilizar Luke para ver el índice y realizar distintas consultas sobre él para demostrar que se ha creado de forma adecuada.

Interfaz de búsqueda

Se debe saber qué y cómo se va a buscar en nuestra interfaz de búsqueda para así saber qué campos de los documentos con los que se trabaja y de qué forma van a ser tenidos en cuenta para elaborar el índice.

Se deben identificar, al menos, los siguientes tipos de campos sobre los documentos de entrada:

- StringField: Texto simple que no se tokeniza. Útil para la búsqueda por facetas y para el filtrado de consultas.
- TextField: Términos que son procesados para la indexación.
- Numérico.
- Facetas (Categorías)

Además, en la aplicación se deben poder realizar:

- Consultas booleanas que involucre a los operadores lógicos OR, AND o NOT.
- Consultas avanzadas (consultas por proximidad, por ejemplo)
- Presentaciones de la información utilizando distintos criterios de ordenación.

Analizador

Hay que identificar qué términos serán utilizados en la búsqueda y cuales no.

Para ello, se debe seleccionar entre los tipos de analizador que tiene implementados Lucene el más adecuado para la aplicación, justificando la decisión. Es posible que haya que utilizar un analizador distinto para cada uno de los campos a indexar.

Facetas

Se deberá realizar la búsqueda por facetas. Para conseguirlo, se deben identificar los campos por los que se podrá clasificar los documentos.

Además, como resultado de la búsqueda podremos tener los resultados agrupados por categorías.

Diseño de la solución

Campo de búsqueda general

Búsqueda de autores

Búsqueda de artículos

Búsqueda en resúmenes

Opciones:

☐ AND

☐ OR

Ordenar por 1	Ordenar por 2	Ordenar por 3	Ordenar por 4	Ordenar por 5
---------------	---------------	---------------	---------------	---------------

Faceta Autores

Faceta Keywords Autores

Faceta Keywords Indexación

Resultados

Nuestra interfaz tendrá un aspecto similar al de la imagen, y, en consecuencia, deberemos implementar la creación del índice de una determinada manera, que describimos a continuación.

En nuestra solución, incluimos todos los campos de los que constan los registros sobre artículos científicos de la base de datos Scopus que hemos utilizado.

De esta forma, podemos devolver toda la información y no nos vemos limitados a la hora de trabajar más adelante con los distintos tipos de consultas.

Para los campos en los que tenemos previsto realizar una búsqueda general o realizar una búsqueda por campo específico utilizamos el analizador EnglishAnalyzer, debido a que tanto el título, como el resumen o los distintos keywords se encuentran en inglés y deben ser tokenizados para facilitar que nuestras búsquedas obtengan buenos resultados (EnglishAnalyzer se ocupa de eliminar palabras vacías, hacer stemming, etc.).

Para el campo EID utilizamos el KeywordAnalyzer para que se indexe la cadena de texto sin dividirla de ningún modo. Esto se hace porque normalmente este campo sólo va a ser usado cuando se conoce el EID de un libro concreto y, por lo tanto, se introduce tal cual en el buscador.

Por último, utilizamos el StandardAnalyzer para el resto de campos, por si es necesario buscar artículos a partir de ellos.

Esta asignación de analizadores se realiza mediante PerFieldanalyzerWrapper en la función cargaAnalyzer.

Hemos seleccionado como facetas los autores de los artículos, las keywords de los autores y las keywords para indexar. Estos valores nos serán muy útiles para clasificar la información que nos devuelvan las búsquedas, ya que tienen un gran poder de agrupación y, además, nos permitirán visualizar la información de forma interesante.

Manual de usuario

Para la implementación del programa que se ocupa de crear el índice se ha utilizado el IDE NetBeans. El mismo se puede utilizar para su ejecución de forma sencilla, para lo cual hace falta especificar como argumento la ruta de la carpeta que contiene los archivos (o la ruta del archivo en el caso de que haya sólo uno).

Para el uso de Luke para la realización de consultas, basta con ejecutarlo y, a continuación, especificar la ruta del índice. Una vez hecho esto, accedemos a la pestaña Search y, una vez ahí, escribimos la expresión deseada además de especificar sobre qué campo queremos realizar la búsqueda y el analizador a utilizar.

En la siguiente imagen podemos observar como ejemplo la búsqueda de la palabra “twitter” en el campo Abstract utilizando el analizador EnglishAnalyzer:

The screenshot shows the Luke - Lucene Index Toolbox (7.1.0) window. The 'Search' tab is active. The search expression is 'twitter'. The analyzer is set to 'org.apache.lucene.analysis.en.EnglishAnalyzer' and the default field is 'Abstract'. The search details show 'Abstract:twitter' and 'Parsed' status. The search results are displayed in a table with columns: #, Score, Doc. Id, Abstract, Author keywords, Authors, Cited by, EID, Index keywords, Link, Source title, and Title. The results show 14 documents, with the first one having a score of 3.6988 and document ID 1154.

#	Score	Doc. Id	Abstract	Author keywords	Authors	Cited by	EID	Index keywords	Link	Source title	Title
0	3,6988	1154	Sentiment	Microblogs;	Giachanou .	13	2-s2.0-8497	Data mining	https://www	ACM Compu	Like it
1	3,6449	112	Twitter mes	DAN2; Feat	Ghiassi M.,	101	2-s2.0-8487	DAN2; Feat	https://www	Expert Syst	Twitter
2	3,6002	1601	Microbloggi	Classifier; O	Shahheidar	9	2-s2.0-8488	Microbloggi	https://www	Proceeding	Twitter
3	3,5844	743	Social medi		Ranco G., A	22	2-s2.0-8494	finance; Int	https://www	PLoS ONE	The eff
4	3,4758	534	Twitter is a	Marijuana; s	Cavazos-Rel	31	2-s2.0-8490	cannabis; a	https://www	Journal of M	Charac
5	3,4529	138	Micro-blogs		Bifet A., Fra	85	2-s2.0-7865	Blogging; D	https://www	Lecture Not	Sentim
6	3,4484	556	The dramat	Box office; F	Arias M., Ari	30	2-s2.0-8489	Box office; I	https://www	ACM Transa	Foreca
7	3,4324	1701	It was not u	Opinion min	Martínez-Cá	8	2-s2.0-8493	Artificial int	https://www	Journal of Ir	Polarit
8	3,4153	1045	To examine		Canvasser I	15	2-s2.0-8492	classificatio	https://www	Journal of E	The us
9	3,4134	2	Twitter is a		Tumasjan A	844	2-s2.0-8489	Content an	https://www	ICWSM 2010	Predict
10	3,4134	302	Twitter as a	Collective ir	Cheong M.,	52	2-s2.0-7404	Blogospher	https://www	Internation	Integra
11	3,4134	400	Twitter sent	Machine lea	Hassan A.,	41	2-s2.0-8489	Multi-class	https://www	Proceeding	Twitter
12	3,4101	965	Twitter is a		Bifet A., Hol	17	2-s2.0-8005	Data strear	https://www	Lecture Not	MOA-TV
13	3,3815	170	Microbloggi		Stieglitz S.,	76	2-s2.0-8485	Social netw	https://www	Proceeding	Politica
14	3,3815	220	The study o	Civilian rec	Cheong M.,	66	2-s2.0-8506	Data visuali	https://www	Information	A micro

Index name: /home/ppm/index/index