

PRÁCTICA III:

...

Implementación de un Sistema de
Recuperación de Información utilizando
Lucene

...

Indexación

26 de octubre de 2022

<i>ÍNDICE</i>	2
---------------	---

Índice

1. Objetivo	3
2. Diseñando el Sistema de Recuperación de Información	3
2.1. Índice en Lucene	4
3. Determinando los campos de indexación de nuestros documentos	6
3.1. Búsqueda por Facetas	7
3.2. Uso de Facetas	9
4. Indexación de documentos	9
4.1. Selección del analizador	12
5. Creación de un Document Lucene	13
6. Segmentos en Lucene	15
7. Ejercicios	16
8. Entrega Práctica	19

1. Objetivo

El objetivo final de esta práctica es que el alumno comprenda todos los procesos que intervienen en el diseño de un Sistema de Recuperación de Información y cómo puede ser implementado utilizando la biblioteca Lucene. Para ello, en este curso utilizaremos la base de datos sobre películas, vista en la práctica anterior. Nuestro objetivo será construir un programa que se pueda ejecutar en línea de comandos y que sea el encargado de generar/actualizar nuestro índice.

En cualquier aplicación de búsqueda podemos distinguir los siguientes pasos, que detallaremos a continuación:

1. Análisis de requisitos y adquisición de datos. (práctica anterior)
2. Procesamiento de los datos, que ya hemos considerado en parte en prácticas anteriores.
3. Indexación y almacenamiento de los mismos. (esta práctica)
4. Búsqueda sobre el índice y presentación de los resultados.

2. Diseñando el Sistema de Recuperación de Información

A la hora de diseñar cualquier aplicación de búsqueda, el primer paso es analizar qué tipo de información se va a buscar y cómo se realizan las búsquedas por parte de un usuario. En esta práctica el objetivo es crear un buscador sobre un conjunto de datos de películas.

En la práctica consideraremos tres tipos de información, no necesariamente excluyente:

- La susceptible de ser tratada como categorías, como el género
- La numérica, como valoraciones, etc.
- La eminentemente textual, que contiene descripción del atributo en lenguaje natural

2 DISEÑANDO EL SISTEMA DE RECUPERACIÓN DE INFORMACIÓN 4

Nosotros nos centraremos en el desarrollo de un sistema de recuperación de información básico capaz de trabajar con estos datos.

En PRADO se puede encontrar un fichero .csv con los datos de entrada. Para poder trabajar con los mismos deberemos ser capaz de extraer los atributos del fichero .csv. Con este fin se recomienda el uso de una librería específica para la lectura de este tipo de información como puede ser openCSV <https://mkyong.com/java/how-to-read-and-parse-csv-file-in-java/> (o cualquier otra que se conozca)

2.1. Índice en Lucene

EL objetivo de esta práctica es construir un índice, pero antes veamos que es un índice Lucene. Conceptualmente es similar a una tabla en una base de datos. La tabla en una base de datos relacional tradicional o en una base de datos NoSQL necesita al menos tener su esquema definido cuando se crea. Hay algunas restricciones claras en la definición de la clave primaria, las columnas y similares. Sin embargo, no existen tales restricciones en un índice Lucene.

Un índice Lucene puede entenderse como un conjunto de documentos. Se puede poner un nuevo documento en el conjunto o se puede sacar uno, pero si se quiere modificar uno de los documentos, primero hay que sacarlo, modificarlo y volver a ponerlo en el conjunto. Puedes meter todo tipo de documentos, y Lucene puede indexar el documento sea cual sea su contenido. Esto nos permitirá incluir dentro del mismo índice documentos de distinta tipología, y realizar la búsqueda sobre los mismos.

Pero, qué es un documento Lucene. Se podría asimilar a una fila en una base de datos relacional. Cuando un documento se inserta en el índice se le asigna un único identificador (DocId). Los documentos se componen de uno o varios campos. Un campo es la unidad más pequeña definible de un índice de datos en Lucene. A cada campo se le asocia un conjunto de valores así como la forma en que los mismos se almacenarán en el índice, que dependerá del tipo de campo (FieldType).

Lucene proporciona muchos tipos diferentes de campos ya creados, aunque podemos personalizarlos y crear nuestros propios tipos de campos.

2 DISEÑANDO EL SISTEMA DE RECUPERACIÓN DE INFORMACIÓN 5

Los atributos que determinan el comportamiento de un campo determinado son:

- **stored:** indica si se guarda el campo. Si es falso, Lucene no almacena el valor del campo y los documentos devueltos en los resultados de la consulta sólo contendrán campos guardados.
- **tokenized:** representa si se tokeniza o no. En Lucene, sólo es necesario tokenizar el campo TextField.
- **termVector:** un vector de términos guarda toda la información relacionada con un término, incluyendo el valor del término, las frecuencias y las posiciones. No se recomienda habilitar el vector de términos para los campos cortos porque toda la información de los términos puede obtenerse tokenizando de nuevo el campo. Sin embargo, se recomienda habilitar el vector de términos para campos más largos o con un alto coste de tokenización.

Hay dos usos principales del vector de términos. El primero es el resaltado de palabras clave (highlighting) y el otro es la comparación de similitudes entre documentos (more-like-this).

- **omitirNormas:** Norms representa la normalización por la longitud del documento. Lucene permite que cada campo de cada documento tenga guardado un factor de normalización, que es un coeficiente que puede afectar al score del documento. Sólo se necesita un byte para guardar Norms, pero se guarda un Norms distinto en cada campo de cada documento y cada dato de Norms se carga en memoria. Así que habilitar las Normas consume espacio de almacenamiento y memoria adicionales. Pero si se desactiva, no podrá utilizarlo en tiempo de consulta.
- **indexOptions:** Lucene proporciona cinco parámetros opcionales para los índices invertidos (NONE, DOCS, DOCS_AND_FREQS, DOCS_AND_FREQS_AND_POSITIONS, DOC_AND_FREQS_AND_POSITIONS_AND_OFFSETS), que se utilizan para seleccionar si el campo necesita ser indexado, y qué contenido indexar.

3 DETERMINANDO LOS CAMPOS DE INDEXACIÓN DE NUESTROS DOCUMENTOS6

- **docValueType:** DocValue es una característica introducida en Lucene 4.0 (asocia los docs a un campo), que mejora enormemente la eficiencia de la clasificación, las facetas y la agregación. DocValues es una estructura de almacenamiento con un esquema fuerte, por lo que todos los campos con DocValues activado deben tener exactamente el mismo tipo. Actualmente Lucene sólo proporciona los cinco tipos de NUMERIC, BINARY, SORTED, SORTED_NUMERIC y SORTED_SET.
- **dimensión:** Lucene soporta la indexación de datos multidimensionales, empleando una indexación especial para optimizar las consultas de datos multidimensionales. El escenario de uso más típico de este tipo de datos es un índice de ubicaciones geográficas. Este es el método de indexación que se utiliza generalmente para los datos de latitud y longitud.

3. Determinando los campos de indexación de nuestros documentos

En gran medida, los campos a considerar depende del uso (tipos de consultas) que un usuario pueda realizar sobre nuestra colección así como la respuesta que daría el sistema.

Los campos concretos dependerá de la aplicación concreta sobre la que estemos trabajando, pero como hemos visto si es necesario un mismo campo podrá utilizarse para dos funciones distintas, por ejemplo como texto sin tokenizar y texto tokenizado.

En cualquier caso, en nuestra aplicación deberemos identificar al menos los siguientes tipos de campos sobre los documentos de entrada:

- **StringField:** Texto simple que se considera literalmente (no se tokeniza), útil para la búsqueda por facetas, filtrado de consultas y también para la presentación de resultados en la interfaz de búsqueda, por ejemplo ID, tags, direcciones webs, nombres de fichero, etc.
- **TextField:** Secuencia de términos que es procesada para la indexación, esto es, pasa por un analyzer pudiendo ser convertida a minúscula, tokenizada,

3 DETERMINANDO LOS CAMPOS DE INDEXACIÓN DE NUESTROS DOCUMENTOS7

estemizada, etc. Como podría ser el título, resumen, etc de un artículo científico.

- Numérico, datos que se expresan mediante este tipo de información, bien sean enteros o reales.
- Facetas (Categorías) que permiten una agrupación lógica de los documentos con la misma faceta, como por ejemplo la revista donde se publicó el trabajo, la fecha de publicación, o el país de origen de los autores.

Además de una consulta por texto libre, en la aplicación se deberá poder realizar como mínimo una consulta booleana que involucre a los operadores lógicos OR, AND o NOT en la misma. Además deberemos proporcionar algún tipo de consulta avanzada como por ejemplo las consultas por proximidad, así como permitir presentar la información utilizando distintos criterios de ordenación (esto es, además de presentar los elementos ordenados por relevancia, debemos de poder presentarlo utilizando un orden distinto).

3.1. Búsqueda por Facetas

Una búsqueda por facetas nos permite acceder a la información refinando la búsqueda de acuerdo a una clasificación por categorías, filtrando los datos teniendo en cuenta las categorías a las que pertenecen. Para ello, es necesario que cada documento pueda ser clasificado a lo largo de múltiples dimensiones (llamadas facetas) como por ejemplo el autor, el idioma, el género o en un sitio de comercio electrónica cada una de las posibles categorías bajo las que podemos clasificar un producto (marca, modelo, características, etc.).

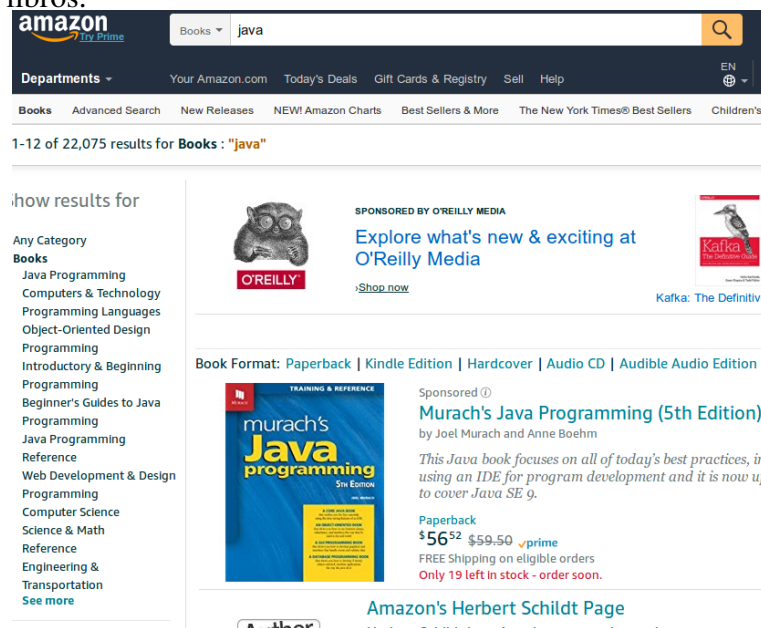
Un atractivo de la mezcla de la búsqueda con el uso de la navegación por facetas es que, ante una consulta, podemos mostrar el número de elementos recuperados en cada una de las categorías. Así, por ejemplo, podemos saber cuantos documentos, de entre los relevantes a la consulta, han sido publicados en el año 2015 o el 2016, o cuántos de ellos han sido protagonizados por un determinado actor. Además, cuando el usuario selecciona una de ellas podemos restringir la búsqueda (drill down) entre los documentos que pertenecen a dicha categoría. Esta información hace fácil la búsqueda de los elementos de interés, ya que

3 DETERMINANDO LOS CAMPOS DE INDEXACIÓN DE NUESTROS DOCUMENTOS8

el usuario puede navegar fácilmente por los resultados, facilitando las siguientes interacciones con el sistema para refinar la búsqueda.

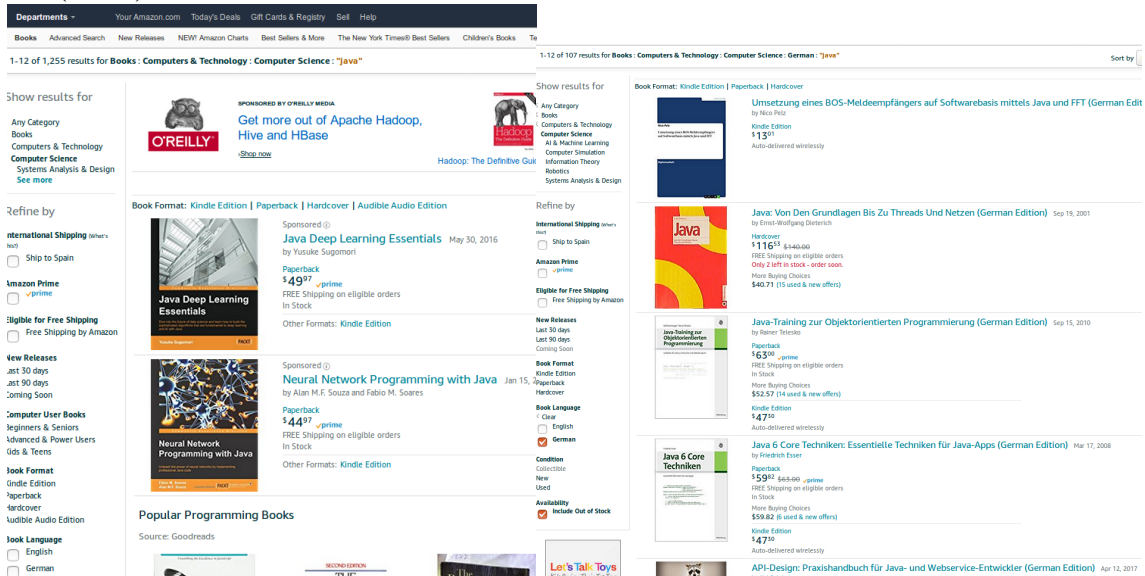
Esta peculiaridad ha hecho que la búsqueda por facetas sea muy común en sitios de comercio electrónico, como por ejemplo Amazon. En la Figura 1 podemos ver cómo ante la consulta “Java” encontramos un total de 22075 libros en el portal de ventas Amazon. A la izquierda de la misma encontramos un frame en el que se permite mostrar los resultados por categorías (aunque Amazon, por motivos internos, ha decidido no mostrar cuántos libros hay en cada una de las categorías). Entre las categorías que considera Amazon encontramos el tipo de libro, lenguaje, autores, formato, etc.

Figura 1: Búsqueda de libros en Amazon, consulta: “Java”. Se encuentran un total de 22.075 libros.



Así, podemos centrar la búsqueda dentro de la categoría Computer-Science (imagen a la izquierda de la Figura 2), encontrando un total de 1255 libros y dentro de ella, podemos de nuevo restringirnos a los libros que han sido editados en Alemán (imagen a la derecha de la Figura 2), encontrando un total de 107 libros.

Figura 2: Consulta: “Java”, restringimos la búsqueda a los libros que se encuadran dentro de la categoría computer-science (izq.) o computer-science -> Alemán (dcha.) Se encuentran un total de 22.075 libros.



3.2. Uso de Facetas

Nuestro sistema final deberá realizar la búsqueda por facetas. Para ello en esta etapa se deberá identificar los campos por los que podrá clasificar los documentos. Así, como resultado de la búsqueda, podremos tener los resultados agrupados por categorías, permitiendo al usuario bucear por ellas en busca de la información de su interés. La documentación más extensa sobre esta parte la veremos en entregas posteriores.

4. Indexación de documentos

Indexar es una de las principales tareas que podemos encontrar en Lucene. La clase que se encarga de la indexación de documentos es `IndexWriter`. Esta clase se encuentra en `lucene-core` y permite añadir, borrar y actualizar documentos Lucene. Un documento Lucene está compuesto por un conjunto de campos: par nombre del campo (string)- contenido del campo (string), como por ejemplo (“autor”, “Miguel de Cervantes”) o (“Título”, “Don Quijote de la Mancha”) o

Cuadro 1: Clases involucradas en la indexación

Clase	Descripción
IndexWriter	Clase esencial que crea/modifica un índice.
Directory	Representa la ubicación del índice.
Analyzer	Es responsable de analizar un texto y obtener los tokens de indexación.
Document	Representa un documento Lucene, esto es, un conjunto de campos (fields) asociados al documento.
Field	La unidad más básica, representa un par clave-valor, donde la clave es el nombre que identifica al campo y el valor es el contenido del documento a indexar.

“(“Cuerpo”, “En un lugar de la Mancha de cuyo...”). Cada uno de estos campos será indexado como parte de un documento, después de pasar por un Analyzer. En el Cuadro 1 podemos ver resumido el conjunto de clases que son usadas frecuentemente en la indexación.

El IndexWriter http://lucene.apache.org/core/9_3_0/core/org/apache/lucene/index/IndexWriter.html toma como entrada dos argumentos:

- Directory: Representa el lugar donde se almacenará el índice http://lucene.apache.org/core/9_3_0/core/org/apache/lucene/store/Directory.html. Podemos encontrar distintas implementaciones, pero para un desarrollo rápido de un prototipo podemos considerar MMapDirectory (el índice se almacena en memoria) o FSDirectory (el índice se almacena en el sistema de ficheros)
- IndexWriterConfig: Almacena la configuración utilizada http://lucene.apache.org/core/9_3_0/core/org/apache/lucene/index/IndexWriterConfig.html

Aunque se recomienda mirar la documentación de las distintas clases, ilustraremos su uso mediante el siguiente ejemplo,

```

1 FSDirectory dir = FSDirectory.open( Paths.get(INDEX_DIR) );
2 IndexWriterConfig config = new IndexWriterConfig( analyzer );
3 config.setOpenMode( IndexWriterConfig.OpenMode.CREATE ).set;
```

```
4      // config.setSimilarity(new LMDirichletSimilarity ());
5
6
7      //Crea un nuevo indice
8      IndexWriter writer = new IndexWriter(dir , config);
9
10     List<string> listaFicheros = .....
11
12     for (String nombre: listaFicheros)
13         Document doc = ObtenerDocDesdeFichero(nombre);
14         writer.addDocument(doc);
15
16     writer.commit();
17     // Ejecuta todos los cambios pendientes en el indice
18     writer.close();
```

Las líneas 1 a 4 son las encargadas de crear las estructuras necesarias para el índice. En concreto el índice se almacenará en disco en la dirección dada por el path. La línea 2 declara un `IndexWriterConfig` `config` con el analizador que se utilizará para todos los campos (si no se indica, utilizará por defecto el `standardAnalyzer`). Para ello, debemos de seleccionar, de entre los tipos que tiene Lucene implementados, el que se considere adecuado para nuestra aplicación. Este es un criterio importante para poder alcanzar los resultados óptimos en la búsqueda. Puede que sea necesario el utilizar un analizador distinto para cada uno de los posibles campos a indexar, en cuyo caso utilizaremos un `PerFieldAnalyzerWrapper` (ver siguiente sección).

La línea 3 se indica que el índice se abre en modo `CREATE` (crea un nuevo índice o sobrescribe uno ya existente), otras posibilidades son `APPEND` (añade documentos a un índice existente) o `CREATE_OR_APPEND` (si el índice existe añade documentos, si no lo crea para permitir la adición de nuevos documentos). Es posible modificar la configuración del índice considerando múltiples setter. Por ejemplo, podremos indicar la función de similitud que se utiliza, por defecto es `BM25` (ver línea 4) o criterios para la mezcla de índices, etc.

Una vez definida la configuración tenemos el `IndexWriter` listo para añadir nuevos documentos. Los cambios se realizarán en memoria y periódicamente se volcarán al `Directory`, que cuando sea necesario realizará la mezcla de distintos

segmentos.

Una vez añadidos los documentos, debemos asegurarnos de llamar a `commit()` para realizar todos los cambios pendientes, línea 16. Finalmente podremos cerrar el índice mediante el comando `close` (línea 18).

4.1. Selección del analizador

El proceso de análisis nos permite identificar qué elementos (términos) serán utilizados en la búsqueda y cuales no (por ejemplo mediante el uso de palabras vacías)

Las operaciones que ejecuta un analizador incluyen: extracción de tokens, supresión de signos de puntuación, acentos o palabras comunes, conversión a minúsculas (normalización), stemización. Para ello, debemos de seleccionar de entre los tipos que tiene Lucene implementados, el que se considere adecuado para nuestra aplicación, justificando nuestra decisión. Este es un criterio importante para poder alcanzar los resultados óptimos en la búsqueda.

En esta práctica será necesario utilizar un analizador distinto al por defecto para algunos de los posibles campos a indexar, para ello podemos utilizar un `PerFieldAnalyzerWrapper` como indica el siguiente ejemplo¹.

```
1 // map field-name to analyzer
2 Map<String , Analyzer> analyzerPerField = new HashMap<String ,
   Analyzer>();
3 analyzerPerField.put("uncampo", new StandardAnalyzer());
4 analyzerPerField.put("otrocampo", new EnglishAnalyzer());
5
6 // create a per-field analyzer wrapper using the Whitespace as
   .. default analyzer ;)
7 PerFieldAnalyzerWrapper analyzer = new PerFieldAnalyzerWrapper(
   new WhitespaceAnalyzer(), analyzerPerField);
```

¹https://lucene.apache.org/core/9_3_0/analysis/common/org/apache/lucene/analysis/miscellaneous/PerFieldAnalyzerWrapper.html

5. Creación de un Document Lucene

Hemos visto que para indexar la información es necesario la creación de un documento, `Document`, Lucene http://lucene.apache.org/core/9_3_0/core/org/apache/lucene/document/Document.html. Un `Document` es la unidad de indexación y búsqueda. Está compuesto por un conjunto de campos `Fields`, cada uno con su nombre y su valor textual. Un `field` puede ser el nombre de un producto, su descripción, su ID, etc. Veremos brevemente cómo se gestionan los `Fields` ya que es la estructura básica de la que está compuesta un `Document`. Información sobre los `Fields` la podemos encontrar en http://lucene.apache.org/core/9_3_0/core/org/apache/lucene/document/Field.html.

Un `Field` tiene tres componentes, el nombre, el valor y el tipo. En nombre hace referencia la nombre del campo (equivaldría al nombre de una columna en una tabla). El valor hacer referencia al contenido del campo (la celda de la tabla) y puede ser texto (`String`, `Reader` o un `TokenStream` ya analizado), binario o numérico. El tipo determina como el campo es tratado, por ejemplo indicar si se almacena (`store`) en el índice, lo que permitirá devolver la información asociada en tiempo de consulta, o si se tokeniza.

Para simplificar un poco la tarea, Lucene dispone de tipos predefinidos

- `TextField`: `Reader` o `String` indexado y tokenizado, sin `term-vector`². Sirve para almacenar el contenido de un documento, sobre el que normalmente realizaremos las búsquedas. Suele pasar por un analizador antes de almacenar la información en el índice.
- `StringField`: Un campo `String` que se indexa como un único token. Por ejemplo, se puede utilizar para ID de un producto, el path donde se encuentra el archivo, etc. Si queremos que la salida de una búsqueda pueda ser ordenada según este campo, se tiene que añadir otro campo del tipo `SortedDocValuesField`.

²En Lucene, un `term-vector` implica que para cada término conocemos: el id del documento, el nombre del campo, el texto en sí, la frecuencia, la posición y los offsets. Esta información podemos hacer cosas interesantes en la búsqueda como conocer dónde emparejan los términos de la consulta (por ejemplo para hacer un `highlighting`)

Para profundizar mas sobre TextField y StringField se puede consultar <https://northcoder.com/post/lucene-fields-and-term-vectors/>

- IntPoint: Entero indexado para búsquedas exactas o por rango. Si queremos devolverlo en consultas, lo debemos de añadir como StoredField
- LongPoint: Long indexado para búsquedas exactas o por rango. Si queremos devolverlo en consultas, lo debemos de añadir como StoredField
- FloatPoint: float indexado para búsquedas exactas o por rango. Si queremos devolverlo en consultas, lo debemos de añadir como StoredField
- DoublePoint: double indexado para búsquedas exactas o por rango. Si queremos devolverlo en consultas, lo debemos de añadir como StoredField.
- SortedDocValuesField: byte[] indexados con el objetivo de permitir la ordenación o el uso de facetas por el campo
- SortedSetDocValuesField: Permite añadir un conjunto de valores, SortedSet, al campo para su uso en facetas, agrupaciones o joinings.
- NumericDocValuesField: Field que almacena a valor long por documento con el fin de utilizarlo para ordenación, facetas o el propio cálculo de scores.
- SortedNumericDocValuedFiled: Añade un conjunto de valores numéricos, SortedSet, al campo
- StoredField: Valores almacenados que solo se utilizan para ser devueltos en las búsquedas.

Los distintos campos de un documento se añaden con el método add.

```
1 Document doc = new Document();
2
3 doc.add(new StringField("isbn", "978-0071809252", Field.Store.YES));
4 doc.add(new TextField("titulo", "Java: A Beginner's Guide, Sixth Edition", Field.Store.YES)); // por defecto no se
    almacena
```

```

5 doc.add(new TextField("contenido", "Fully updated for Java
    Platform , Standard Edition 8 (Java SE 8), Java ....."));
6 doc.add(new IntPoint("size", 148));
7 doc.add(new StoredField("size",148));
8 doc.add(new SortedSetDocValuesField("format", new BytesRef("
    paperback")));
9 doc.add(new SortedSetDocValuesField("format", new BytesRef("
    kindle")));
10
11 Date date = new SimpleDateFormat("yyyy-MM-dd'T'HH:mm:ss'Z').
    parse(...);
12 doc.add(new LongPoint("Date", date.getTime()));

```

Como podemos imaginar, para cada tipo tenemos un comportamiento específico, aunque nosotros podremos crear nuestro propio tipo de campo.

```

1
2 FieldType authorType = new FieldType();
3 authorType.setIndexOptions(IndexOptions.DOCS_AND_FREQS);
4 // authorType.setIndexOptions(IndexOptions.
5     DOCS_AND_FREQS_AND_POSITIONS_AND_OFFSETS);
6 authorType.setStored(true);
7 authorType.setOmitNorms(true);
8 authorType.setTokenized(false);
9 authorType.setStoreTermVectors(true);
10
11 doc.add(new Field("author", "Arnaud Cogoluegues", authorType))
12     ;
13 doc.add(new Field("author", "Thierry Templier", authorType));
14 doc.add(new Field("author", "Gary Gregory", authorType));

```

6. Segmentos en Lucene

Veamos con un poco mas de detalle como Lucene gestiona internamente un índice.

Un índice se compone de uno o varios subíndices. Un subíndice se denomina segmento. El diseño conceptual de los segmentos en Lucene es similar al de LSM, pero hay algunas diferencias. Heredan las ventajas de escritura de datos de un

LSM (log-structured merge tree), pero sólo proporcionan consultas en tiempo casi real y no en tiempo real.

Cuando Lucene escribe datos, primero lo hace en un buffer en memoria (similar a MemTable en LSM). Cuando los datos del búfer alcanzan una determinada cantidad, se vacían para convertirse en un segmento. Cada segmento tiene su propio índice independiente y se puede buscar de forma independiente, pero los datos nunca pueden ser modificados. Este esquema evita las escrituras aleatorias. Los datos se escriben en lotes (Batch) se añaden (Append) alcanzando un alto rendimiento. Los documentos escritos en el Segmento no pueden ser modificados, pero pueden ser borrados. El método de borrado no modifica el archivo en su ubicación interna original, pero el DocID del documento que se va a borrar se guarda en otro archivo para garantizar que el archivo de datos no pueda ser modificado.

Las consultas al índice necesitan consultar varios segmentos y fusionar los resultados, así como tener en cuenta los documentos eliminados. Para optimizar las consultas, Lucene tiene una política para fusionar múltiples segmentos y, en este sentido, es similar a la fusión de SSTable de LSM.

Hasta que la información no esta volcada al índice, los datos se almacenan en la memoria y no se pueden buscar. Esta es otra razón por la que se dice que Lucene proporciona consultas en tiempo casi real y no en tiempo real.

7. Ejercicios

1. Basándonos en el código anterior, implementar un pequeño programa que nos permite añadir varios documentos a un índice Lucene, podemos seguir el siguiente esquema.

```
1
2  import org.apache.lucene.analysis.Analyzer;
3  import org.apache.lucene.analysis.standard.
      StandardAnalyzer;
4  ....
5
6  public class IndiceSimple {
7
8      String indexPath = "./index";
```



```
9      String docPath = "./DataSet";
10
11      boolean create = true;
12
13      private IndexWriter writer;
14
15      public static void main(String[] args) {
16          // Analizador a utilizar
17          Analyzer analyzer = new StandardAnalyzer();
18          // Medida de Similitud (modelo de recuperacion) por
              defecto BM25
19          Similarity similarity = new ClassicSimilarity();
20          // Llamados al constructor con los parametros
21          IndiceSimple baseline = new IndiceSimple( ... );
22
23          // Creamos el indice
24          baseline.configurarIndice(analyzer, similarity);
25          // Insertar los documentos
26          baseline.indexarDocumentos();
27          // Cerramos el indice
28          baseline.close();
29
30
31      }
32
33      public void configurarIndice(Analyzer analyzer, Similarity
          similarity) throws IOException {
34
35
36          IndexWriterConfig iwc = new IndexWriterConfig(
              analyzer);
37          iwc.setSimilarity(similarity);
38          // Crear un nuevo indice cada vez que se ejecute
39          iwc.setOpenMode(IndexWriterConfig.OpenMode.CREATE);
40          // Para insertar documentos a un indice existente
41          // iwc.setOpenMode(IndexWriterConfig.OpenMode.
              CREATE_OR_APPEND);
42
43      }
```

```
44
45      // Localizacion del indice
46      Directory dir= FSDirectory.open(Paths.get(indexPath)
47      );
48
49      // Creamos el indice
50      writer = new IndexWriter(dir, iwc);
51
52  }
53
54  public void indexarDocumentos() {
55
56      // Para cada uno de los documentos a insertar
57      for (elementos d : docPath) {
58
59          //leemos el documento sobre un string
60          String cadena = leerDocumento( d );
61          // creamos el documento Lucene
62          Document doc = new Document();
63
64          //parseamos la cadena (si es necesario)
65          Integer start,end; // Posiciones inicio,fin
66
67          // Obtener campo Entero de cadena
68          Integer start = ... // Posicion de inicio del campo;
69          Integer end = ... // Posicion fin del campo;
70          String aux = cadena.substring(start, end);
71          Integer valor = Integer.decode(aux);
72
73          // Almacenamos en el campo en el documento Lucene
74          doc.add(new IntPoint("ID", valor));
75          doc.add(new StoredField("ID", valor));
76
77          // Obtener campo texto de cadena
78          start = ... // Posicion de inicio del campo;
79          end = ... // Posicion fin del campo;
80          String cuerpo = cadena.substring(start,end);
81          // Almacenamos en el campo en el documento Lucene
```

```
82         doc.add(new TextField("Body", cuerpo, Field.Store.YES));
83
84         //Obtenemos los siguientes campos
85
86         // .....
87
88         // Insertamos el documento Lucene en el indice
89         writer.addDocument(doc);
90         // Si lo que queremos es actualizar el documento
91         // writer.updateDocument(new Term("ID", valor.toString()), doc);
92     }
93
94 }
95
96
97 public void close() {
98     try {
99         writer.commit();
100        writer.close();
101    } catch (IOException e) {
102        System.out.println("Error closing the index.");
103    }
104 }
105
106 }
```

2. Utilizar Luke para ver el índice y realizar distintas consultas sobre el mismo.

8. Entrega Práctica

La fecha para la entrega de la práctica es el viernes 11 de noviembre de 2022. El ejecutable deberá permitir o bien crear un índice desde el principio o añadir posibles nuevas películas a una colección ya creada (asumiendo el mismo formato).

Los grupos de dos alumnos deberán permitir la inclusión en el índice de otros datos que siguiesen un esquema diferente, como por ejemplo los que nos podemos

descargar de Kaggle <https://www.kaggle.com/datasets/cryptexcode/mpst-movie-plot-synopses-with-tags>, también en formato csv o de CMU <http://www.cs.cmu.edu/~ark/personas/>, aunque en este caso tendrán que fusionar la información disponible en distintos ficheros.