

Sistema de Recuperación de Información



Jose Manuel Rodríguez Calvo
Juan Miguel Vilchez Rodríguez

Índice

1. Indexación
2. Búsqueda
3. Facetas
4. Bibliografía

Indexación

Esta parte es una de las más importantes de nuestro proyecto y en la que más tiempo hemos invertido debido a que una buena indexación dará lugar a unos buenos resultados en las búsquedas. Los tres documentos con los que empezamos son **Questions**, **Answers** y **tags**, todos los utilizamos para ofrecer unos mejores resultados y entre ellos tienen un formato diferente que hace que se aplique diferentes métodos para almacenar la información, sin embargo usamos una indexación lucene para los tres. Los tratamientos para los diferentes archivos son los siguientes:

77434,14008,2008-09-16T21:40:29Z,171,How to access the last value in a vector? , " <p> Suppose I have a vector that is nested in a data frame one or two levels . Is there a quick and dirty way to access the last value , without using the <code>length()</code> function? Something a la PERL's <code>>\$ # </code> specialvar?</p> <p>So I would like something like : </p> <pre><code> dat \$vec1 \$vec2 [\$ #] </code> </pre> <p> instead of </p> <pre><code> dat \$vec1 \$vec2 [length(dat \$vec1 \$vec2)] </code> </pre> "</code>	All: StoredField Id_Q: StringField UserId: StringField DateLong: LongPoint → para ser consultado Date: StoredField ScoreInt: IntPoint → para ser consultado Score: StoredField Title: TextField Body: TextField Code: TextField
--	--

Además de indexarlos, hemos usado un **PerFieldAnalyzerWrapper** para analizar según el campo ya que no tiene que recibir el mismo tratamiento el título que el código. Tanto para el campo **Title** como para el **body** hemos usado **EnglishAnalyzer** ya que todos están en Inglés, mientras que para el campo **Code** hemos desarrollado un analizador propio para que sea más exacto con la necesidad de información.

Para la indexación de las respuestas, usamos un método algo diferente aunque el trasfondo del funcionamiento es muy similar, llegando así a reutilizar código implementado en el apartado anterior. A continuación muestro un ejemplo del tratamiento de los datos

<pre> 9336125,1216960,2012-02- 17T22:07:41Z,9336071,2,FALSE,"<p>Don't know if this exactly answers what you are asking, but:</p> <pre><code>vector &lt;- res\$par &gt; vector [1] 0.9771973 1.0218072 &gt; vector[1] [1] 0.9771973 </code></pre> <p>There may be an accessor method that is more 'proper.'</p> </pre>	<pre> All: StoredField Id_A: StringField UserId: StringField DateLong: LongPoint → para ser consultado Date: StoredField Id_parent: StringField ScoreInt: IntPoint → para ser consultado Score: StoredField Aproved: StringField Body: TextField Code: TextField </pre>
---	---

El tratamiento de los datos, en concreto sobre el analizador que le aplicamos a cada campo es el mismo que preguntas. la única diferencia son los campos: **Aproved** e **Id_parent** ; a los que no le hemos aplicado ningún analizador debido a que son tipos de datos unitarios donde solo son datos simples que son buscados de forma booleana.

Para el tratamiento de los **tags**, son considerados como parejas de id de la pregunta y el tag, por tanto el tratamiento de este documento es sencillo, donde incluso solo aplicamos a tag el analizador **WhiteSpaceAnalyzer** porque el otro campo será buscado de manera unitaria y sin alterar ese valor:

77434,vector	<pre> Id_T : StringField tag: StringField </pre>
--------------	--

A parte de estos campos hemos tenido que añadir otro campo más para poder conseguir un mejor tratamiento de los datos que es el campo **Type** que es un *StringField* en el que almacenamos el tipo del documento sobre el que estamos tratando, así cuando estamos ante una pregunta el valor de **Type** será **"Question"**, el de Answer con su correspondiente **"Answer"** y tag con **"tag"**.

Con todo este proceso funcional, el problema que intentamos solucionar se vuelve más sencillo y cómodo de implementar posteriores pasos. Sobre el rendimiento de este proceso que se realiza antes de mostrarle los resultados a los usuarios, por tanto el tiempo de esta parte no es relevante para nuestro problema. Para una mejor monitorización de nuestro sistema, con un intel core i7 de 2.80 GHz y una memoria RAM de 32GB este proceso tarda alrededor de tres minutos y ocupa un total de 815 MB.

Para la ejecución de esta parte '**core.java**' será muy importante la carpeta **index**, que se creará en el mismo repositorio y si no está creada esa carpeta la creará el programa, en cambio si ya está creada, el programa interpreta que ya no hay que hacer el índice y no procederá a ejecutarlo.

Búsqueda

Este es otro de los puntos muy importantes para obtener buenos resultados a la necesidad de información, en este proceso podemos distinguir dos funciones importantes, la **búsqueda normal** (Junto con la búsqueda avanzada) y la **búsqueda por facetas**.

La búsqueda normal hace acopio de una función especializada que busca el *query* deseado en todos los posibles campos de los documentos almacenados devolviéndolos ordenados por **Score**. Dado que *Stack Overflow* tiene su propio **Score** reflejado en cada documento, hemos extraído manualmente este para mezclarlo con el **Score** de lucene y así reflejar lo que creemos que es un orden por puntuación más realista.

La búsqueda avanzada hace una consulta más específica, a la que se le pasa una lista con las *queries* deseadas por cada campo, por ejemplo, podríamos buscar documentos que tengan el *query* deseado en el código, o en el título, o podríamos buscar fácilmente por usuarios concretos. Esto se hace mediante una combinación de BooleanClauses.

Aquí tenemos ejemplos de funciones que usamos y resultados de las búsquedas.

```
public static Map searchin_process(String consultation, String Facets, String FacetType, String [] listaFacetas) throws Exception{
    Set<Document> solutions_All = new HashSet<Document>();
    Set<String> solutions_Tags = new HashSet<String>();
    //test that consultation is not a file
    File af = new File(consultation);
    String document="";
    String line="";
    if (af.isFile()){
        BufferedReader br = new BufferedReader(new FileReader(consultation));
        while ((line = br.readLine()) != null) {
            document=document+" "+line;
        }
        body_titlequery(document,solutions_All);
        codequery(document,solutions_All);
    }
    else if(Integer.parseInt(Facets) == 1){
        facetsQuery(consultation,solutions_All, Integer.parseInt(FacetType));
    }
    else if(Integer.parseInt(Facets) == 2){
        facetsQuerySearch(consultation, listaFacetas,solutions_All, Integer.parseInt(FacetType));
    }
    else{
        //test if the query is an ID of user, an ID of question/answer or a tag
        if(consultation.indexOf(" ")!=-1){
            unitariquery(consultation,solutions_All,solutions_Tags);
        }
        else{
            body_titlequery(consultation,solutions_All);
            codequery(consultation,solutions_All);
        }
    }
}
```

Este es el método principal de las búsquedas, al que se le pasa la consulta, un string ("1" o "2") que diferencia una búsqueda normal de una por facetas, String FacetType, que nos dice el tipo de faceta que estamos buscando, y una lista de facetas para una vez elegida poder buscar con esos nuevos filtros. Como podemos ver, tenemos una comprobación para ver si la consulta es un único término (en cuyo caso le pasamos un filtro para comprobar si se está buscando

un usuario o una ID) y luego otro tipo de consulta para si tenemos más de un término (Que busca en el título, en el cuerpo y en el código)

En este ejemplo vemos como se construyen las clausulas booleanas para la búsqueda :

```
QueryParser parserID_Q =new QueryParser("Id_Q",new EnglishAnalyzer());
Query q1= parserID_Q.parse(query);
QueryParser parserID_A =new QueryParser("Id_A",new EnglishAnalyzer());
Query q2= parserID_A.parse(query);
QueryParser parserID_T =new QueryParser("Id_T",new EnglishAnalyzer());
Query q3= parserID_T.parse(query);
QueryParser parserUser =new QueryParser("User",new EnglishAnalyzer());
Query q4= parserUser.parse(query);
QueryParser parsertag =new QueryParser("tag",new EnglishAnalyzer());
Query q5= parsertag.parse(query);

BooleanClause bc1 = new BooleanClause(q1, BooleanClause.Occur.SHOULD);
BooleanClause bc2 = new BooleanClause(q2, BooleanClause.Occur.SHOULD);
BooleanClause bc3 = new BooleanClause(q3, BooleanClause.Occur.SHOULD);
BooleanClause bc4 = new BooleanClause(q4, BooleanClause.Occur.SHOULD);
BooleanClause bc5 = new BooleanClause(q5, BooleanClause.Occur.SHOULD);

BooleanQuery.Builder bqbuilder = new BooleanQuery.Builder();
bqbuilder.add(bc1);
bqbuilder.add(bc2);
bqbuilder.add(bc3);
bqbuilder.add(bc4);
bqbuilder.add(bc5);

BooleanQuery bq = bqbuilder.build();

TopDocs docs=searcher.search(bq,20);
```

Aquí unos pantallazos de algunos resultados obtenidos en las búsquedas con la aplicación.

The screenshot shows a web application titled "BUSCADOR". On the left, there is a search interface with a "Buscar" button, a "Vector" input field, and a "Buscar Facetas" button. Below these are two radio buttons for "Fecha" and "Autores", followed by a list of 10 search results, each with a radio button and a snippet of text. At the bottom left is a button labeled "Buscar con las nuevas opciones". On the right, there is a section titled "BUSQUEDA AVANZADA" with input fields for "Título:", "Cuerpo:", "Usuario:", and "Código:", and a "Buscar" button. The main content area displays a large snippet of code, which appears to be a Java program snippet related to vector operations and iteration.

BUSCADOR

+vector +program +java

☐ Fecha ☐ Autores

- ☐
- ☐
- ☐
- ☐
- ☐
- ☐
- ☐
- ☐
- ☐
- ☐

```

33211810.5296362,2015-10-19T10:11:55Z,0,"Java loading error:""character vector expected""",<p>I am using 64 bit windows 7 &amp; 64 bit R 3.2.0. I have down
lockquote>
<p>>Error : onLoad failed in loadNamespace() for 'Jlava', details:<br>
at: dimName(this$RuntimeLib), error: a character vector argument
expected in addition. Warning message: package 'Jlava' was built under
version 3.2.2. Error: package or namespace load failed for 'Jlava'</p>
lockquote>
<p>I searched stackoverflow for answers, and did the following-</p>
<p>->1) Added the <code>jvm.dll</code> file to my PATH variable, i.e. in the following directory:</p>
re<code>C:\Program Files\Java\jdk1.8.0_25\bin
code>-</pre>
<p>->2) I am using 64 bit R 3.2.0 &amp; checked if the JVM version was 64 bit or not, typing the following on the command line:</p>
re<code>java -d64 -version
ode>-</pre>
<p>->It gave the following output-</p>
re<code>C:\Users\wy2&gt;.java -d64 -version
a version "1.8.0_25"
(TM) SE Runtime Environment (build 1.8.0_25-b18)
a HotSpot(TM) 64-Bit Server VM (build 25.25-b02, mixed mode)
ode>-</pre>
<p>->So at least the JVM matches with R.</p>
<p>->So what more needs to be done to make <code>JJava</code> work? I also checked my PATH variable and it seems okay.</p>

=====
re<code>vector &lt;i>c</i>: c("dog","ant","eagle","ant","eagle","parrot")
dog""ant""eagle""ant""eagle""parrot"
ipping &lt;i>df</i>-data.frame(key=c("dog","cat","elephant","ant","parrot","eagle"),value=c("mammal","mammal","mammal","insect","bird","bird"))
by value
og mammal
at mammal
ephant mammal
nt insect

```

BUSQUEDA AVANZADA

Título:

Cuerpo:

Usuario:

Código:

[illegible]

Facetas

La inclusión de facetas en nuestro buscador nos permite hacer búsquedas más específicas y mejorar la interacción del usuario con el proyecto. Nosotros hemos implementado dos facetas diferentes, que son, la **fecha** de los documentos y los **autores** de los mismos, permitiendo al usuario sesgar la búsqueda con ellos.

Primero debemos indexar las facetas deseadas, lo hacemos junto con la indexación del resto de campos de la siguiente manera :

```
FSDirectory taxoDir = FSDirectory.open(Paths.get(FacetsIndexPath));
FacetsConfig fconfig = new FacetsConfig();

DirectoryTaxonomyWriter taxoWriter = new DirectoryTaxonomyWriter(taxoDir);
fconfig.setHierarchical("Publish_Date", true);
fconfig.setMultiValued("Author", true);
```

Aquí abrimos el índice de facetas y generamos su configuración, creando la faceta de carácter jerárquico **Publish_Date** y el campo multivalued **Author**.

```
doc.add(new FacetField("Author", separate[1]));
doc.add(new FacetField("Publish_Date", String.valueOf(year)));

writer.addDocument(fconfig.build(taxoWriter, doc));
```

Una vez creados, para cada documento indexamos su valor de faceta, para el autor la ID del creador del documento y para la fecha, nos quedamos con el año.

En la parte de búsqueda de nuestro programa, creamos un apartado específico para las facetas. Que son buscadas con la siguiente función :

```
DrillDownQuery ddq = new DrillDownQuery(fconfig, bq);
for (String prueba : Facet){
    System.out.println("Autor : " + prueba);
}

if(faceta==1){
    for (String facetaAutor : Facet){
        if(facetaAutor!=null){
            ddq.add("Author", facetaAutor);
        }
    }
}
else{
    for (String facetaFecha : Facet){
        if(facetaFecha!=null){
            ddq.add("Publish_Date", facetaFecha);
        }
    }
}

FacetsCollector fc = new FacetsCollector();
TopDocs tdc = FacetsCollector.search(searcher, ddq, 10, fc);

Facets facetas = new FastTaxonomyFacetCounts(taxoReader, fconfig, fc);
List<FacetResult> TodasDims = facetas.getAllDims(20);

for (ScoreDoc sd : tdc.scoreDocs){
    Document d=searcher.doc(sd.doc);
    solutions_All.add(d);
}
```


Usaremos a continuación un pantallazo de la búsqueda por facetas para explicar a nivel de usuario como funcionan.

BUSCADOR

Buscar

Buscar Facetas

☒ Fecha ☐ Autores

Fecha : 2016, valor (#n)-> 10442 ☒ 1

Fecha : 2015, valor (#n)-> 9780 ☒ 2

Fecha : 2017, valor (#n)-> 8090 ☒ 3

Fecha : 2014, valor (#n)-> 8071 ☐ 4

Fecha : 2013, valor (#n)-> 6383 ☐ 5

Fecha : 2012, valor (#n)-> 3688 ☐ 6

Fecha : 2011, valor (#n)-> 1958 ☐ 7

Fecha : 2010, valor (#n)-> 776 ☐ 8

Fecha : 2009, valor (#n)-> 189 ☐ 9

Fecha : 2008, valor (#n)-> 6 ☐ 10

Buscar con las nuevas opciones

BUSQUEDA AVANZADA

Título :

Cuerpo :

Usuario :

Código :

Buscar

```
020599913279625 30904455,4890647,2015-06-18T00:05:33Z,2,Vectorize thinking,"<p>I'm having a vectorization problem. Say I'm having a vector<code>x</code> &lt;- c(0
#####
>For example,</p>

what I want to do is iterate through this set of data and then place the data contained in these rows into a vector of vectors for processing later. The data contained in
re<code>1-&gt;1-&gt;1
#####
&gt;2
code</pre>
>So row 1 would contain only data from 1 type of ID, then the next row in the vector would be a vector of another type of ID. How would I go about doing this in R? In C+
>Is this even the right way to be approaching this problem? Is there a better way to do what I'm trying to do?</p>
#####
re<code>&gt;foo &lt;- seq(from=1, to=5, by=1)
t_bar &lt;- seq(from=6, to=10, by=1)
t_baz &lt;- seq(from=11, to=15, by=1)
t_vecs &lt;- c(foo,bar,baz)
t_for(v in vecs) {print(v)}
1
2
3
4
5
6
7
8
9
10
#####
```

Primero buscamos en la pestaña “**Buscar facetas**” el query deseado, en nuestro caso es “**vector**”, nos saldrán los resultados para dicho query usando una búsqueda normal, pero a la izquierda, tendremos un recuento (Gracias a las facetas) de en este caso, las fechas en las que más documentos se han creado. Por ejemplo queremos los documentos que contengan vector, de los tres años en los que más documentos fueron creados (2016, 2015, 2017), y ahora relanzamos la búsqueda con los nuevos parámetros.

BUSCADOR

Buscar

Buscar Facetas

☒ Fecha ☐ Autores

Fecha : 2016, valor (#n)-> 10442 ☒ 1

Fecha : 2015, valor (#n)-> 9780 ☒ 2

Fecha : 2017, valor (#n)-> 8090 ☒ 3

Fecha : 2014, valor (#n)-> 8071 ☐ 4

Fecha : 2013, valor (#n)-> 6383 ☐ 5

Fecha : 2012, valor (#n)-> 3688 ☐ 6

Fecha : 2011, valor (#n)-> 1958 ☐ 7

Fecha : 2010, valor (#n)-> 776 ☐ 8

Fecha : 2009, valor (#n)-> 189 ☐ 9

Fecha : 2008, valor (#n)-> 6 ☐ 10

Buscar con las nuevas opciones

BUSQUEDA AVANZADA

Título :

Cuerpo :

Usuario :

Código :

Buscar

```
02009999919435 38807349,NA,2016-08-06T18:10:52,3,vector of punctuation,"<p>For digits I can write a vector like this: </p>
re<code>digits &lt;- c("0","1","2","3","4","5","6","7","8","9")
code</pre>
>How can I get an analogous vector of punctuation marks?</p>
#####
>Still, in his example he demonstrates 'only' how to append one data piece at a time. I am now fighting with the idea to fill a vector with vectors - not scalars.</p>
>Imagine I have a vector with a length of 100</p>
re<code>vector &lt;- numeric(length=100)
code</pre>
>and a smaller vector that would fit 10 times into the first vector</p>
re<code>vec &lt;- seq(1,10,1)
code</pre>
>How would I have to construct a loop that adds the smaller vector to the large vector without using c() or append ?</p>
>EDIT: This example is simplified - vec does not always consist of the same sequence but is generated within a for loop and should be added to vector.</p>
#####
re<code>## Some data
l &lt;- data.frame(p=1:10, b=1:20, c=21:30)
Vectorize with mapgl, seems to work
l<-function(i, j, dat) list(dat[i, j])
lply(l, list(1,2, 3,4), list(1,2, 2,3), MoreArgs = list(dat=dat))
[]
a b
1 11
#####
```

Ahora los resultados obtenidos son de estas tres fechas. En el caso de los usuarios es la misma funcionalidad, nos aparecen los usuarios referentes a ese tema que más documentos han publicado, y podemos filtrar las búsquedas para ver sus respuestas o preguntas.

Bibliografía

Material e información esencial para la indexación, búsqueda y las facetas :

https://decsai.ugr.es/index.php?p=material_alumno

Página principal de Lucene 7.10 API, donde hemos encontrado las bibliotecas necesarias para el funcionamiento de la práctica y métodos ya desarrollados en lucene que han facilitado la misma :

https://lucene.apache.org/core/7_1_0/core/overview-summary.html

Stack Overflow ha sido usado para resolver dudas sobre sintaxis o problemas concretos :

<https://es.stackoverflow.com/>

Video tutoriales para aprender a manejar JFrame y poder desarrollar la interfaz gráfica de la aplicación :

<https://www.youtube.com/watch?v=HC6jljJ1484>

<https://www.youtube.com/watch?v=xCpkFcZ6Yfw>

https://www.youtube.com/watch?v=C2zP_UN43x4

<https://www.youtube.com/watch?v=2dAWfs1NZns>

<https://netbeans.org/kb/docs/java/gui-functionality.html>

<https://www.youtube.com/watch?v=5ynxP-P-nG0>

Instalación de NeatBeans y puesta a punto para la práctica :

<https://netbeans.org/community/releases/82/install.html>

Páginas de información adicional sobre lucene :

<http://www.lucenetutorial.com/lucene-in-5-minutes.html>

https://lucene.apache.org/core/7_1_0/index.html