



# Proyecto RI

## Recuperación de la Información

Cándido Alonso Barro  
Andrés Picazo Mesa

**CURSO 2024/25**



# Índice general

<b>1. Crawler</b>	<b>5</b>
1.1. Descarga documentos . . . . .	6
<b>2. Indexación</b>	<b>7</b>
2.1. Preprocesado del texto . . . . .	7
2.2. Creación del índice invertido y cálculo de longitud . . . . .	7
2.2.1. Cálculo de TF . . . . .	8
2.2.2. Cálculo de IDF y longitud . . . . .	8
<b>3. Búsqueda</b>	<b>9</b>
3.1. La interfaz . . . . .	9
3.2. Ranking de documentos . . . . .	10
3.3. Cómo diferencia entre AND, OR o ninguno . . . . .	11
3.4. Cálculo de documentos relevantes AND . . . . .	11
3.5. Cálculo de puntuajes de documentos relevantes . . . . .	12
3.6. Mostrar resultados . . . . .	12
<b>4. Puntos extra</b>	<b>13</b>
4.1. Stemming . . . . .	13
4.2. Operadores AND y OR en la búsqueda . . . . .	13
<b>5. Instalación y manual del usuario</b>	<b>15</b>
5.1. Instalacion . . . . .	15
5.2. Manual del usuario . . . . .	15
<b>6. Anexo</b>	<b>17</b>
6.1. Preprocesado.java . . . . .	17
6.2. Indexacion.java . . . . .	19
6.3. CorpusCrawler.java . . . . .	21
6.4. Buscador.java . . . . .	23



# Capítulo 1

## Crawler

Este código define un crawler llamado `CorpusCrawler` que utiliza la librería `Jsoup` para descargar y analizar enlaces desde una página HTML alojada en un repositorio remoto. Los archivos encontrados en la página se descargan y organizan en el sistema local manteniendo su estructura relativa.

### Constantes principales

- `BASE_URL`: URL base del repositorio remoto donde se encuentran los archivos.
- `START_URL`: URL inicial que actúa como punto de entrada para el análisis.
- `OUTPUT_DIRECTORY`: Directorio local donde se almacenarán los archivos descargados.

### Metodo principal

Aquí podemos observar lo importante del crawler, es cuando descarga y analiza la página inicial. Primero se conecta a la URL proporcionada anteriormente y descarga su contenido, extrae todos los elementos `<a>` que contengan `href` y ya dentro del bucle descarga cada archivo con su URL adecuada.

Listing 1.1: Estructuras de datos

```
1 System.out.println("Descargando índice...");  
2 Document document = Jsoup.connect(START_URL).get();
```

```
3 Elements links = document.select("a[href]");
4
5 for (Element link : links) {
6     String relativePath = link.attr("href");
7     String fileUrl = BASE_URL + relativePath;
8     downloadFile(fileUrl, relativePath);
9 }
```

### 1.1. Descarga documentos

La función `downloadFile` descarga un archivo desde una URL proporcionada y lo almacena localmente. Para lograrlo, establece una conexión con la URL del archivo y obtiene el flujo de entrada desde el cual leer los datos. Antes de guardar el archivo, verifica y crea los subdirectorios que hagan falta según la ruta. Luego, escribe los datos en un archivo local usando un buffer. Si la descarga y el almacenamiento salen bien, muestra un mensaje en la consola diciendo que el archivo se ha descargado correctamente. En caso de que ocurra un error muestra un mensaje de error.

# Capítulo 2

## Indexación

### 2.1. Preprocesado del texto

El preprocesamiento del texto se realiza a través de una clase dedicada (`Preprocesado.java`), diseñada para mejorar la modularidad y legibilidad del código. Esta clase recibe un objeto de tipo `String` y, mediante la función `procesar(String)`, aplica una serie de operaciones para normalizar el texto. Estas operaciones incluyen la conversión a minúsculas, la eliminación de signos de puntuación, números, espacios innecesarios, guiones, barras invertidas y palabras vacías (stopwords).

La lista de palabras vacías se carga desde un archivo externo (`stopwords-en.txt`) utilizando una estructura `HashSet` para garantizar una búsqueda eficiente.

La clase `Preprocesado.java` se puede ver en: Anexo-6.1

### 2.2. Creación del índice invertido y cálculo de longitud

Para almacenar el índice invertido en memoria durante el proceso de indexación, se ha utilizado un `HashMap` en el que las claves corresponden a cada uno de los términos lematizados del corpus. Cada clave está asociada a un valor que consiste en una `Tupla`, implementada en `Tupla.java`. El primer elemento de la tupla representa el IDF del término, mientras que el segundo es otro `HashMap` que contiene los nombres de los documentos en los que aparece dicho término (como clave) y el peso correspondiente de este término en cada documento (como valor).

Esta estructura de datos permite una mayor eficiencia tanto en las operaciones de

búsqueda como en las de inserción, a cambio de un mayor uso de memoria durante la ejecución.

La longitud de cada documento se almacena en otro Hashmap, esta vez más simple, donde la clave son los nombres de los documentos y su valor asociado su longitud.

Listing 2.1: Estructuras de datos

```
1 private static Map<String, Tupla<Double, Map<String, Double  
    >>> indice_invertido = new HashMap<String, Tupla<Double,  
    Map<String, Double>>>();  
2 private static Map<String, Double> longitud = new HashMap<>();
```

### 2.2.1 Cálculo de TF

Como se indica en el documento de ayuda del proyecto, el cálculo del TF-IDF se realiza en dos pasos.

Primero se calcula el tf de cada término en cada documento aprovechando el bucle que lee todos los ficheros del corpus (6.2). Para obtener la frecuencia de cada término dentro de un documento se emplea la función `dividir_en_terminos()` (6.3) que emplea una estructura `HashMap<String, Integer>` auxiliar, que se vaciará antes de pasar al siguiente documento.

Por último, con la frecuencia obtenida se calcula el TF (6.4) y se empieza a rellenar la estructura que empleará el índice invertido.

### 2.2.2 Cálculo de IDF y longitud

En este paso se calcula el IDF de cada término, la longitud de cada documento y el peso de cada término en cada documento. Se recorre cada término almacenado en el Hashmap y se obtienen los valores necesarios para el cálculo del IDF, número de documentos del corpus (obtenido al recorrer los documentos al principio) y el número de documentos en los que aparece dicho término (empleo la función `size()` del Hashmap que devuelve el número de claves).

Ahora se calcula el peso de término y se realizan los cálculos parciales para la longitud. Por último, se aplica la raíz cuadrada a cada valor de la longitud. La función encargada de esto se puede ver en: Anexo-6.6.



# Capítulo 3

## Búsqueda

### 3.1. La interfaz

El programa comienza con una interfaz basada en texto que permite al usuario ingresar consultas de búsqueda. La interfaz guía al usuario mediante mensajes iniciales y acepta consultas que pueden incluir operadores como AND u OR. También cuenta con un comando para salir del programa escribiendo “salir”. A continuación, se presenta un fragmento del código relevante:

```
1 System.out.print(">");
2 String query = scanner.nextLine().trim();
3
4 // Salir si el usuario escribe "salir"
5 if (query.equalsIgnoreCase("salir")) {
6     System.out.println("Saliendo del buscador.");
7     break;
8 }
9 // Procesar la consulta
10 query = preprocesado.procesar(query);
11
12 // Realizar la búsqueda y ranking
13 Map<String, Double> rankResultado = rankDocumentos(query);
14
15 // Mostrar los resultados
16 mostrarResultados(rankResultado);
17 }
18 scanner.close();
```

## 3.2. Ranking de documentos

El ranking de documentos se realiza a partir de una consulta procesada. Se calculan puntuaciones basadas en la frecuencia del término (TF) y el factor inverso de frecuencia de documentos (IDF). Los documentos se ordenan en función de su relevancia y se devuelven los mejores resultados. El cálculo sería tal que:

$puntuacion = tfDocumento * idfDocumento;$

```
1 private static Map<String, Double> rankDocumentos(String
2     consulta) {
3     Stemmer stemmer = new Stemmer();
4     String[] terminos = consulta.split(" ");
5     Map<String, Double> puntuacionesDocumentos = new HashMap
6         <>();
7
8     for (String termino : terminos) {
9         stemmer.add(termino.toCharArray(), termino.length());
10        stemmer.stem();
11        termino = stemmer.toString();
12
13        List<DocumentoPeso> pesosDocumentos = indiceInvertido
14            .getOrDefault(termino, Collections.emptyList());
15        for (DocumentoPeso docPeso : pesosDocumentos) {
16            double puntuacion = docPeso.tf * docPeso.idf;
17            puntuacionesDocumentos.merge(docPeso.
18                nombreDocumento, puntuacion, Double::sum);
19        }
20    }
21
22    return puntuacionesDocumentos.entrySet().stream()
23        .sorted((a, b) -> Double.compare(b.getValue(), a.
24            getValue()))
25        .limit(10)
26        .collect(Collectors.toMap(
27            Map.Entry::getKey,
28            Map.Entry::getValue,
29            (e1, e2) -> e1,
30            LinkedHashMap::new
31        ));
32 }
```

### 3.3. Cómo diferencia entre AND, OR o ninguno

La consulta ingresada se analiza para identificar si contiene los operadores AND o OR. En este caso se puede observar que los analizo para cuando están en minúsculas y es porque antes ya han recibido un trato de preprocesamiento. Dependiendo de esto, se procesan los términos de manera distinta:

- **AND:** Los documentos relevantes deben contener todos los términos especificados.
- **OR:** Los documentos relevantes pueden contener cualquiera de los términos especificados.
- **Sin operadores:** Se tomara como si fuera un termino unico.

El siguiente fragmento de código ilustra esta lógica:

```
1 boolean esConsultaAnd = consulta.contains("and");
2 boolean esConsultaOr = consulta.contains("or");
3 String[] terminos;
4
5 if (esConsultaAnd) {
6     terminos = consulta.split("\\sand\\s");
7 } else if (esConsultaOr) {
8     terminos = consulta.split("\\sor\\s");
9 } else {
10     terminos = new String[]{consulta.trim()};
11 }
```

### 3.4. Cálculo de documentos relevantes AND

Cuando la consulta tiene el operador AND, se realiza una intersección de los documentos asociados a cada término. Así aseguramos que solo los documentos que contienen los términos sean considerados relevantes.

```
1 if (esConsultaAnd) {
2     if (documentosRelevantes.isEmpty()) {
3         documentosRelevantes = new HashSet<>(
4             documentosParaTermino);
5     } else {
6         documentosRelevantes.retainAll(documentosParaTermino)
7     }
8 }
```

```
6     }  
7 }
```

### 3.5. Cálculo de puntuajes de documentos relevantes

Para cada término de la consulta, se calcula un puntaje que se basa en la fórmula. Los puntajes se suman y se van acumulando para cada documento relevante.

```
1 for (DocumentoPeso docPeso : pesosDocumentos) {  
2     if (documentosRelevantes.contains(docPeso.nombreDocumento  
3         )) {  
4         double puntuacion = docPeso.tf * docPeso.idf;  
5         puntuacionesDocumentos.merge(docPeso.  
6             nombreDocumento, puntuacion, Double::sum);  
7     }  
8 }
```

### 3.6. Mostrar resultados

Los resultados se presentan ordenados por su puntuacion en orden descendente. Si no se encuentran documentos relevantes, se muestra.

```
1 private static void mostrarResultados(Map<String, Double>  
2     resultadosRankeados) {  
3     if (resultadosRankeados.isEmpty()) {  
4         System.out.println("No se encontraron documentos  
5             relevantes para la consulta.");  
6     } else {  
7         System.out.println("Documentos encontrados:");  
8         resultadosRankeados.forEach((doc, puntuacion) ->  
9             System.out.printf("- %s (Puntuacion: %.4f)%n", doc,  
10                 puntuacion));  
11     }  
12 }
```

# Capítulo 4

## Puntos extra

### 4.1. Stemming

Se ha aplicado el algoritmo de Porter para el stemming empleando una clase Java encontrada en GitHub. El repositorio se puede acceder mediante el siguiente enlace: <https://gist.github.com/ldclakmal/667d8ecb620a0cce7d3dedae80a2c013>.

El algoritmo de Porter funciona principalmente eliminando sufijos típicos del inglés, como pueden ser '-ing', '-ed', '-ly'. Aunque no siempre se pueda obtener la raíz exactamente correcta empleando el algoritmo, consideramos que es una solución óptima para nuestro sistema RI debido a su simplicidad y velocidad, aparte de ser un algoritmo muy probado y extendido.

### 4.2. Operadores AND y OR en la búsqueda

Como ya vimos anteriormente en la parte de búsqueda, el AND y el OR se obtienen viendo si dichas palabras están dentro de la consulta. En el caso de que estén, si es un AND vemos qué documentos comparten en común y son esos los que elegimos, y en el caso de que sea un OR, escogemos todos. Para el cálculo de la puntuación de los documentos que comparten ambos, se hace mediante la operación:

$$\text{puntuación} = tf_{Documento1} \cdot idf_{Documento1} + tf_{Documento2} \cdot idf_{Documento2}$$



# Capítulo 5

## Instalación y manual del usuario

### 5.1. Instalacion

Para ejecutar el proyecto, lo que necesitamos hacer es asegurarnos de tener instalado *Java* en nuestro ordenador. Puedes verificarlo ejecutando el siguiente comando en la terminal:

```
java --version
```

Si no te aparece la versión, tendrás que instalarte *Java*. Una vez que tengas *Java* instalado, descarga y descomprime la carpeta del proyecto que te hemos dado. Dentro de esta carpeta, encontrarás todos los archivos. Para compilar y ejecutar el proyecto, solo necesitas ejecutar el archivo `ejecucion.bat`, que está incluido en la carpeta. Este archivo es un script que automatiza todo el proceso de compilación y ejecución del proyecto.

El corpus también debe encontrarse en la carpeta raíz del proyecto. De lo contrario, será necesario modificar la variable `corpus_path` de `Indexacion.java` y sustituirla por la dirección deseada.

### 5.2. Manual del usuario

Al ejecutar el archivo `ejecucionIndexacion.bat`, se ejecutará tanto el crawler como la parte de indexación, y seguidamente tendremos que ejecutar el archivo `ejecucionBuscador.bat`, que ejecutará el buscador.

## Recuperación de la Información

---

Cuando el programa esté en funcionamiento, se te pedirá que ingreses las palabras o términos que deseas buscar. Para hacerlo, solo tendrás que escribir la palabra clave y presionar Enter. El programa realizará la búsqueda en los datos procesados y te devolverá los resultados relevantes.

Si en algún momento deseas finalizar el programa, simplemente escribe la palabra 'salir' y presiona Enter. El programa se cerrará de inmediato.



# Capítulo 6

## Anexo

A continuación se presentan algunos de los códigos desarrollados:

### 6.1. Preprocesado.java

Listing 6.1: Clase Preprocesado.java

```
1 public class preprocesado {
2
3     private static Set<String> STOPWORDS = new HashSet<>();
4
5     static {
6         try {
7             STOPWORDS = new HashSet<>(Files.readAllLines(
8                 Paths.get("./utility/stopwords-en.txt")));
9         } catch (IOException e) {
10             e.printStackTrace();
11         }
12     }
13
14     public static String procesar(String cad) {
15         cad = minusculas(cad);
16         cad = eliminar_signos(cad);
17         cad = eliminar_barra_invertida(cad);
18         cad = eliminar_numeros(cad);
19         cad = eliminar_guiones(cad);
20         cad = eliminar_espacios(cad);
21         cad = eliminar_stopwords(cad);
22     }
23 }
```

```
21
22     return cad;
23 }
24
25 private static String eliminar_signos(String cad) {
26     Pattern pat = Pattern.compile("[!\"#$%&'()
27         *+,./:;<=>?\\@\\[\\]\\^_`{|}~]");
28     Matcher mat = pat.matcher(cad);
29     return mat.replaceAll("_");
30 }
31
32 private static String eliminar_barra_invertida(String cad
33 ) {
34     return cad.replace("\\", "_");
35 }
36
37 private static String eliminar_guiones(String cad) {
38     return cad.replaceAll("\\s-+\\s", "_");
39 }
40
41 private static String eliminar_numeros(String cad) {
42     return cad.replaceAll("\\b\\d*\\b", "_");
43 }
44
45 private static String minusculas(String cad) {
46     return cad.toLowerCase();
47 }
48
49 private static String eliminar_espacios(String cad) {
50     return cad.replaceAll("\\s+", "_");
51 }
52
53 private static String eliminar_stopwords(String cad) {
54     ArrayList<String> aux = Stream.of(cad.split("\\s"))
55         .collect(Collectors.toCollection(ArrayList<
56             String>::new));
57     aux.removeAll(STOPWORDS);
58     return aux.stream().collect(Collectors.joining("_"));
59 }
```

## 6.2. Indexacion.java

Listing 6.2: Bucle que lee los documentos

```

1 File dir = new File(corpus_path);
2     if (dir.exists() && dir.isDirectory()) {
3         File[] documentos = dir.listFiles();
4         if (documentos != null) {
5             for (File documento : documentos) {
6                 try {
7                     String contenido = new String(Files.
8                         readAllBytes(Paths.get(documento.
9                             getPath())));
10                    contenido = preprocesado.procesar(
11                        contenido);
12                    dividir_en_terminos(contenido);
13                    calcular_tf(documento.getName());
14                    terminos_map.clear();
15                } catch (IOException e) {
16                    e.printStackTrace();
17                }
18                // Aumento el numero de documentos
19                N++;
20            }
21        }
22    }

```

Listing 6.3: Division del texto en términos y conteo de la frecuencia de estos

```

1 private static void dividir_en_terminos(String texto) {
2     // Divido el texto en palabras
3     String[] terminos = texto.split("\\s+");
4     // Sin esto se crea siempre al principio un registro
5     // vacio
6     ArrayList<String> listaTerminos = new ArrayList<>(Arrays.
7         asList(terminos));
8     listaTerminos.removeIf(String::isEmpty);
9     terminos = listaTerminos.toArray(new String[0]);
10    // Recorro todos los terminos
11    for (String termino : terminos) {
12        //Aplico el algoritmo de stemming
13        Stemmer stemmer = new Stemmer();
14        char[] termArray = termino.toCharArray();

```

```
13     stemmer.add(termArray, termArray.length);
14     stemmer.stem();
15     termino = stemmer.toString();
16     // Si no esta en el map
17     if (terminos_map.get(termino) == null)
18         // Inicializo el valor a 1
19         terminos_map.put(termino, 1);
20     // Si esta
21     else {
22         // Sumo 1 al valor
23         terminos_map.put(termino, terminos_map.get(
24             termino) + 1);
25     }
26 }
```

Listing 6.4: Cálculo del TF

```
1 private static void calcular_tf(String name) {
2     // Recorro todos los terminos del documento actual
3     for (Map.Entry<String, Integer> entry : terminos_map.
4         entrySet()) {
5         // Obtengo el termino
6         String termino = entry.getKey();
7         // Obtengo su frecuencia
8         Integer frecuencia = entry.getValue();
9         // Calculo el tf
10        Double tf = 1 + Math.log(frecuencia) / Math.log(2);
11        // Si no esta el termino en el mapa tf-idf lo
12        inicializo vacio
13        if (!indice_invertido.containsKey(termino)) {
14            Map<String, Double> mapaInterno = new HashMap<>()
15            ;
16            Tupla<Double, Map<String, Double>> nuevaTupla =
17            new Tupla<>(0.0, mapaInterno);
18            // Inserta la nueva tupla en el mapa con la clave
19            proporcionada
20            indice_invertido.put(termino, nuevaTupla);
21        }
22        // Recupero la tupla del termino actual
23        Tupla<Double, Map<String, Double>> tupla_actual =
24        indice_invertido.get(termino);
25        // Recupero el map donde guardo el documento y su tf
26        Map<String, Double> map_actual = tupla_actual.second;
```

```
21         // Guardo el id y el tf
22         map_actual.put(name, tf);
23     }
24 }
```

Listing 6.5: Cálculo del IDF y longitud

```
1 private static void calcular_idf_y_longitud() {
2     for (Map.Entry<String, Tupla<Double, Map<String, Double
3         >>> entry : indice_invertido.entrySet()) {
4         Tupla<Double, Map<String, Double>> tupla_actual =
5             entry.getValue();
6         Integer n = tupla_actual.second.size();
7         Double a = (double) N / n;
8         Double idf = Math.log(a) / Math.log(2);
9         tupla_actual.first = idf;
10        for (Map.Entry<String, Double> doc : tupla_actual.
11            second.entrySet()) {
12            String docName = doc.getKey();
13            Double peso = doc.getValue() * idf;
14            doc.setValue(peso);
15            if (longitud.get(docName) != null)
16                longitud.put(docName, longitud.get(docName) +
17                    peso * peso);
18            else
19                longitud.put(docName, peso * peso);
20        }
21    }
22    for (Map.Entry<String, Double> doc : longitud.entrySet())
23        doc.setValue(Math.sqrt(doc.getValue()));
24 }
```

## 6.3. CorpusCrawler.java

Listing 6.6: Crawler

```
1 public class CorpusCrawler {
2     private static final String BASE_URL = "https://raw.
3         githubusercontent.com/PdedP/RECINF-Project/refs/heads/
4         main/";
5     private static final String START_URL = BASE_URL + "index
6         .html"; // Pagina inicial
```

```
4     private static final String OUTPUT_DIRECTORY = ".";
5
6     public static void main(String[] args) {
7         try {
8             File directory = new File(OUTPUT_DIRECTORY);
9             if (!directory.exists()) {
10                 directory.mkdir();
11             }
12
13             // Descargar y analizar la pagina inicial
14             System.out.println("Descargando índice...");
15             Document document = Jsoup.connect(START_URL).get
16                 ();
17             Elements links = document.select("a[href]");
18
19             for (Element link : links) {
20                 String relativePath = link.attr("href");
21                 String fileUrl = BASE_URL + relativePath;
22                 downloadFile(fileUrl, relativePath);
23             }
24
25             System.out.println("Descarga completada. Los
26                 archivos están en el directorio: " +
27                 OUTPUT_DIRECTORY);
28         } catch (IOException e) {
29             System.err.println("Error al procesar el índice:
30                 " + e.getMessage());
31         }
32     }
33
34     private static void downloadFile(String fileUrl, String
35         relativePath) {
36         try (InputStream in = new URL(fileUrl).openStream())
37         {
38             // Crea los subdirectorios
39             File outputFile = new File(OUTPUT_DIRECTORY,
40                 relativePath);
41             outputFile.getParentFile().mkdirs();
42
43             try (OutputStream out = new FileOutputStream(
44                 outputFile)) {
45                 byte[] buffer = new byte[1024];
46                 int bytesRead;
```

```
39         while ((bytesRead = in.read(buffer)) != -1) {
40             out.write(buffer, 0, bytesRead);
41         }
42         System.out.println("Archivo descargado: " +
43             relativePath);
44     }
45 } catch (IOException e) {
46     System.err.println("Error al descargar el archivo
47         : " + fileUrl + "- " + e.getMessage());
48 }
```

## 6.4. Buscador.java

Listing 6.7: Buscador de documentos

```
1 import utility.Stemmer;
2 import utility.preprocesado;
3
4 import java.io.*;
5 import java.util.*;
6 import java.util.stream.Collectors;
7
8 public class Buscador {
9
10     // Indice invertido construido desde el archivo
11     private static final Map<String, List<DocumentoPeso>>
12         indiceInvertido = new HashMap<>();
13
14     public static void main(String[] args) {
15         // Construir el indice desde el archivo en la carpeta
16         "utility"
17         try {
18             cargarIndexArchivo("utility/indice_invertido.dat"
19                 );
20         } catch (IOException e) {
21             System.err.println("Error al cargar el indice
22                 invertido: " + e.getMessage());
23             return; // Termina el programa si no se puede
24                 cargar el archivo
25         }
26     }
27 }
```

```
20     }
21
22     Scanner scanner = new Scanner(System.in);
23     System.out.println("Bienvenido al buscador de
24         documentos.");
25     System.out.println("Introduce una consulta (puedes
26         usar operadores AND/OR):");
27
28     while (true) {
29         System.out.print(">");
30         String query = scanner.nextLine().trim();
31
32         // Salir si el usuario escribe "salir"
33         if (query.equalsIgnoreCase("salir")) {
34             System.out.println("Saliendo del buscador.");
35             break;
36         }
37
38         // Procesar la consulta
39         query = preprocesado.procesar(query);
40
41         // Realizar la búsqueda y ranking
42         Map<String, Double> rankResultado =
43             rankDocumentos(query);
44
45         // Mostrar los resultados
46         mostrarResultados(rankResultado);
47     }
48
49     scanner.close();
50 }
51
52 private static void cargarIndexArchivo(String fileName)
53     throws IOException {
54     try (BufferedReader reader = new BufferedReader(new
55         FileReader(fileName))) {
56         String line;
57         while ((line = reader.readLine()) != null) {
58             // Parsear cada línea del archivo
59             String[] parts = line.split(";");
60             if (parts.length < 3) continue; // Saltar
61                 líneas mal formateadas
```



```
57         String word = parts[0].toLowerCase(); //
           Primera parte: palabra clave
58         double idf = Double.parseDouble(parts[1]); //
           Segunda parte: IDF
59         List<DocumentoPeso> documents = new ArrayList
           <>();
60
61         // Procesar documentos-tf asociados
62         for (int i = 2; i < parts.length; i++) {
63             String[] docParts = parts[i].split("-");
64             if (docParts.length > 1) {
65                 String nombreDocumento = docParts[0];
66                 double tf = Double.parseDouble(
67                     docParts[1]);
68                 documents.add(new DocumentoPeso(
69                     nombreDocumento, tf, idf));
70             }
71         }
72         // Agregar al indice invertido
73         indiceInvertido.put(word, documents);
74     }
75 }
76
77 private static Map<String, Double> rankDocumentos(String
78 query) {
79     // Crear una instancia del stemmer
80     Stemmer stemmer = new Stemmer();
81
82     // Determinar el tipo de consulta
83     boolean esAndQuery = query.contains("and");
84     boolean esOrQuery = query.contains("or");
85
86     // Si la consulta tiene AND o OR, dividirla en
87     terminos
88     String[] terminos;
89     if (esAndQuery) {
90         terminos = query.split("\\s*and\\s*"); // Se
           maneja AND con espacios opcionales
91     } else if (esOrQuery) {
92         terminos = query.split("\\s*or\\s*"); // Se
           maneja OR con espacios opcionales
```

```
91     } else {
92         // Si no contiene AND ni OR, tratamos toda la
           consulta como un unico termino
93         terminos = new String[]{query.trim()}; // Se
           toma como un unico termino
94     }
95
96     Set<String> documentoRelevante = new HashSet<>();
97     Map<String, Double> puntuacionDocumento = new HashMap
           <>(); // Para acumular los puntajes
98
99     for (String term : terminos) {
100         term = term.trim(); // Limpiar posibles espacios
           extras
101
102         if (term.isEmpty()) continue; // Asegurarse de
           que no estamos procesando terminos vac os
103
104         // Aplicar stemming al termino
105         stemmer.add(term.toCharArray(), term.length());
106         stemmer.stem();
107         term = stemmer.toString(); // Obtener la raiz
           del termino
108
109         List<DocumentoPeso> pesoDocumento =
           indiceInvertido.getOrDefault(term, Collections.
           emptyList());
110         Set<String> documentParaTermino = pesoDocumento.
           stream()
111             .map(dw -> dw.nombreDocumento)
112             .collect(Collectors.toSet());
113
114         if (esAndQuery) {
115             // Si es una consulta con AND, hacemos la
           interseccion de documentos
116             if (documentoRelevante.isEmpty()) {
117                 documentoRelevante = new HashSet<>()
           {
118                     addAll(documentParaTermino);
119                 }
120             } else {
121                 documentoRelevante.retainAll(
           documentParaTermino);
122             }
123         } else if (esOrQuery) {
124             documentoRelevante.addAll(
           documentParaTermino);
125         }
126     }
127
128     // Calcular la puntuacion para cada documento
129     for (String doc : documentoRelevante) {
130         double suma = 0.0;
131         for (String term : terminos) {
132             suma += puntuacionDocumento.getOrDefault(
           term, 0.0);
133         }
134         puntuacionDocumento.put(doc, suma);
135     }
136
137     // Ordenar los documentos por puntuacion
138     List<String> documentosOrdenados = new ArrayList<>();
139     for (Map.Entry<String, Double> entry : puntuacionDocumento.
           entrySet()) {
140         documentosOrdenados.add(entry.getKey());
141     }
142
143     // Devolver los documentos ordenados
144     return documentosOrdenados;
145 }
```

```
122         // Si es una consulta con OR, hacemos la
123         union de documentos
124         documentoRelevante.addAll(documentParaTermino
125         );
126     } else {
127         // Implicito OR si no hay operadores
128         documentoRelevante.addAll(documentParaTermino
129         );
130     }
131 }
132
133 // Si no hay documentos relevantes, retornar un mapa
134 vacio
135 if (documentoRelevante.isEmpty()) {
136     return Collections.emptyMap();
137 }
138
139 // Calcular puntajes para los documentos relevantes
140 for (String term : terminos) {
141     term = term.trim(); // Limpiar espacios
142
143     if (term.isEmpty()) continue; // Saltar si el
144     termino esta vacio
145
146     // Aplicar stemming nuevamente (si es necesario)
147     antes de calcular puntajes
148     stemmer.add(term.toCharArray(), term.length());
149     stemmer.stem();
150     term = stemmer.toString();
151
152     List<DocumentoPeso> pesoDocumento =
153         indiceInvertido.getOrDefault(term, Collections.
154         emptyList());
155     for (DocumentoPeso docWeight : pesoDocumento) {
156         if (documentoRelevante.contains(docWeight.
157         nombreDocumento)) {
158             // Calcular el puntaje para cada
159             documento
160             double score = docWeight.tf * docWeight.
161             idf;
162             // Sumar puntaje al documento sin
163             duplicar
```

```
152         puntuacionDocumento.merge(docWeight.  
153             nombreDocumento, score, Double::sum);  
154     }  
155 }  
156  
157 // Ordenar los documentos por puntaje en orden  
158 // descendente  
159 return puntuacionDocumento.entrySet()  
160     .stream()  
161     .sorted((a, b) -> Double.compare(b.getValue()  
162         , a.getValue()))  
163     .limit(10) // Mostrar solo los 10 mejores  
164     resultados  
165     .collect(Collectors.toMap(  
166         Map.Entry::getKey,  
167         Map.Entry::getValue,  
168         (e1, e2) -> e1,  
169         LinkedHashMap::new  
170     ));  
171 }  
172  
173 private static double calculaIDF(String term) {  
174     // Obtener los documentos donde aparece el termino  
175     List<DocumentoPeso> pesoDocumento = indiceInvertido.  
176         getOrDefault(term, Collections.emptyList());  
177     if (pesoDocumento.isEmpty()) return 0.0;  
178  
179     // Tomar el IDF del primer documento asociado  
180     return pesoDocumento.get(0).idf;  
181 }  
182  
183 private static void mostrarResultados(Map<String, Double  
184 > rankResultado) {  
185     if (rankResultado.isEmpty()) {  
186         System.out.println("No se encontraron documentos  
187             relevantes para la consulta.");  
188     } else {  
189         System.out.println("Documentos encontrados:");  
190         rankResultado.forEach((doc, score) -> System.out.  
191             printf("- %s (Score: %.4f)%n", doc, score));  
192     }  
193 }
```

```
187
188 // Clase auxiliar para representar un documento y su peso
189 // Clase auxiliar para representar un documento, su TF y
    el IDF del termino
190 private static class DocumentoPeso {
191     String nombreDocumento;
192     double tf; // Frecuencia del termino en el documento
193     double idf; // IDF del termino
194
195     DocumentoPeso(String nombreDocumento, double tf,
        double idf) {
196         this.nombreDocumento = nombreDocumento;
197         this.tf = tf;
198         this.idf = idf;
199     }
200 }
201
202 }
```