

# Count Word

Andres Ricardo Caceres Ortiz 506222067

**Abstract**—Se ha propuesto desarrollar algoritmos para el procesamiento de texto. En este taller se propondrán dos ejemplos de algoritmos para la búsqueda y enumeración de las palabras más recurrentes en un documento previamente definido. No solo se identificarán las palabras más repetidas, sino también se mostrará la ubicación precisa de cada palabra en el documento. También se calculará la complejidad de cada algoritmo, en este caso se hallarán la complejidad espacial y temporal de cada algoritmo, de este modo garantizar una eficiencia óptima tanto en términos de tiempo como de espacio.

## I. INTRODUCTION

La manipulación y el análisis de grandes cantidades de texto se han vuelto esenciales en una variedad de aplicaciones, desde motores de búsqueda en la web hasta la minería de datos y la inteligencia artificial. Es por ello que se ha propuesto la creación de dos algoritmos, capaces de separar cada palabra y almacenarla en un arreglo, así mismo este deberá mostrar la ubicación exacta y las veces que se repite. Antes de desarrollar el taller se deben tener en cuenta algunas definiciones como:

HashMap: es una estructura de datos en programación que representa una colección de pares clave-valor, donde cada clave está asociada a un valor.

Arreglo o Array: Es un tipo de dato que almacena datos, dichos datos son del mismo tipo (String, int etc), todo esto de manera estructurada.

Array List: Es una clase similar a los Array, con la diferencia de que para un ArrayList no es necesario declarar su tamaño.

## II. ALGORITMO 1: MEMORIZACIÓN PARA LA CLASIFICACIÓN DE FRECUENCIA DE PALABRAS

Calcular las frecuencias de palabras puede ser computacionalmente costoso, especialmente cuando se trabaja con conjuntos de datos voluminosos. Para abordar este problema, podemos emplear un algoritmo que utiliza la memorización, una técnica de caché que almacena resultados previamente calculados para evitar cálculos redundantes.

El primer paso consiste en crear un HashMap para almacenar pares de palabra-frecuencia. A continuación, recorreremos cada documento, lo dividimos en palabras y actualizamos las frecuencias de las palabras en el HashMap. Después de procesar todos los documentos, tenemos un HashMap que contiene pares de palabra-frecuencia. Para identificar las palabras más repetidas, ordenamos el HashMap por valores (frecuencias) en orden descendente. Este enfoque garantiza que las palabras más frecuentes aparezcan en la parte superior de la clasificación.

La caché de memorización almacena el HashMap de palabra-frecuencia calculado, evitando la necesidad de recalcularlo para los mismos datos de entrada en el futuro. Esta optimización no solo ahorra recursos computacionales, sino que también acelera el proceso de identificación de términos clave dentro del texto. A continuación el desarrollo del algoritmo:

```
1
2 import java.util.*;
3
4 public class DocumentSearch {
5
6     private static Map<String, List<
7         Integer>> invertedIndex = new
8         HashMap<>();
9     private static Map<String, List<
10         Integer>> cache = new HashMap<>();
11     public static void buildInvertedIndex
12         (String[] documents) {
13         for (int i = 0; i < documents.
14             length; i++) {
15             String document = documents[i]
16                 .split("\\s+");
17             for (String word : words) {
18                 word = word.toLowerCase();
19                 if (!invertedIndex.
20                     containsKey(word)) {
21                     invertedIndex.put(
22                         word, new
23                         ArrayList<>());
24                 }
25                 invertedIndex.get(word).
26                     add(i);
27             }
28         }
29     }
30
31     public static List<Integer>
32         searchWord(String word) {
33         if (cache.containsKey(word)) {
34             return cache.get(word);
35         }
36
37         word = word.toLowerCase();
38         List<Integer> result =
```

```

29         invertedIndex.getOrDefault(
30             word, new ArrayList<>());
31
32     cache.put(word, result);
33
34     return result;
35 }
36
37 public static void main(String[] args)
38 {
39     String[] myDocuments = {
40         "La programaci n en Python
41         es clave para el trabajo
42         con datos",
43         "Los programadores en Java
44         tienen un alto inters en
45         pasar a Python",
46         "La optimizaci n de
47         algoritmos es fundamental
48         en el desarrollo de
49         software",
50         "...";
51     };
52
53     buildInvertedIndex(myDocuments);
54
55     String searchWord = "desarrollo";
56     List<Integer> documentIndices =
57         searchWord(searchWord);
58
59     System.out.print(searchWord + ":
60     ");
61     for (int i = 0; i <
62         documentIndices.size(); i++) {
63         System.out.print(
64             documentIndices.get(i)+1);
65         if (i < documentIndices.size
66             () - 1) {
67             System.out.print(", ");
68         }
69     }
70     System.out.println("]");
71 }
72 }

```

#### A. Complejidad Temporal y Espacial

- Complejidad Temporal: En el peor caso, este algoritmo debe recorrer todos los documentos en la colecci3n, lo que resulta en una complejidad de  $O(n)$ , donde  $n$  es el n3mero de documentos. La b3squeda en cada documento se realiza en tiempo lineal en funci3n de la longitud del documento y la longitud de la palabra objetivo, por lo que podemos considerar la b3squeda dentro de un

documento como  $O(m)$ , donde  $m$  es la longitud m3xima del documento o la palabra objetivo m3s larga.

- Complejidad Espacial: El espacio adicional utilizado por este algoritmo es principalmente para almacenar la lista de documentos que contienen la palabra objetivo. En el peor caso, podr3a ser  $O(n)$ , donde  $n$  es el n3mero de documentos que contienen la palabra objetivo.

### III. ALGORITMO 2: 3NDICE INVERTIDO PARA RECUPERACI3N DE DOCUMENTOS

El segundo algoritmo que exploraremos realiza la b3squeda de documentos que contienen una palabra espec3fica en un conjunto de documentos. La cach3 juega un papel crucial en este algoritmo. Una vez construido el 3ndice invertido, podemos almacenar en cach3 los resultados, asegurando que las b3squedas posteriores de las mismas palabras clave sean extremadamente r3pidas. Cuando se solicita una b3squeda, primero verificamos la cach3 y, si se encuentra el resultado, lo devolvemos directamente.

```

1 import java.util.*;
2 import java.util.stream.Collectors;
3
4 public class WordFrequencyRanking {
5     private static Map<String[], Map<
6         String, Integer>> cache = new
7         HashMap<>();
8
9     public static Map<String, Integer>
10         getWordFrequency(String[]
11             documents) {
12         if (cache.containsKey(documents))
13             {
14                 return cache.get(documents);
15             }
16         Map<String, Integer>
17             wordFrequency = new HashMap
18             <>();
19         for (String document : documents)
20             {
21                 String[] words = document.
22                     split("\\s+");
23                 for (String word : words) {
24                     word = word.toLowerCase();
25                     ;
26                     wordFrequency.put(word,
27                         wordFrequency.
28                         getOrDefault(word, 0)
29                         + 1);
30                 }
31             }
32         Map<String, Integer>
33             sortedWordFrequency =
34             wordFrequency.entrySet().
35                 stream().
36                 sorted(Map.Entry.<String
37                     , Integer>

```

```

22         comparingByValue().
           reversed())
           .collect(Collectors.toMap(
           (Map.Entry::getKey,
           Map.Entry::getValue, (
           e1, e2) -> e1,
           LinkedHashMap::new));
23
24
25     cache.put(documents,
           sortedWordFrequency);
26
27     return sortedWordFrequency;
28 }
29
30 public static void main(String[] args
31 ) {
32     String[] myDocuments = {
33         "La programacin en Python
           es clave para el trabajo
           con datos",
34         "Los programadores en Java
           tienen un alto inters en
           pasar a Python",
35         "La optimizacin de
           algoritmos es fundamental
           en el desarrollo de
           software",
36         "...
37     };
38     Map<String, Integer>
           wordFrequency =
           getWordFrequency(myDocuments);
39
40
41     int rank = 1;
42     for (Map.Entry<String, Integer>
           entry : wordFrequency.entrySet
           ()) {
43         System.out.println(entry.
           getKey() + " " + entry.
           getValue());
44         rank++;
45         if (rank > 10) {
46             break;
47         }
48     }
49 }
50 }

```

#### A. Complejidad Temporal

Construcción del Índice Invertido (buildInvertedIndex):  $O(N * M)$ , donde  $N$  es el número de documentos y  $M$  es el número

promedio de palabras en un documento. Búsqueda de Palabras (searchWord):  $O(1)$  en promedio debido al uso de la caché. Sin embargo, en el peor caso, donde la palabra no está en la caché y debe buscarse en el Índice Invertido, la complejidad sería  $O(N)$ , donde  $N$  es el número promedio de documentos donde aparece la palabra.

#### B. Complejidad Espacial

Espacio para el Índice Invertido:  $O(W * D)$ , donde  $W$  es el número total de palabras únicas en todos los documentos y  $D$  es el número de documentos. Espacio para la Caché:  $O(W)$ , donde  $W$  es el número total de palabras únicas.

### IV. CONCLUSIONES

En este artículo, hemos explorado dos algoritmos para buscar palabras en una colección de documentos y hemos discutido su complejidad temporal y espacial. La elección del algoritmo adecuado depende de varios factores, como el tamaño de la colección de documentos y las operaciones de búsqueda previstas. La Búsqueda Lineal es simple pero puede volverse ineficiente para colecciones grandes, mientras que la Búsqueda con Mapa es más eficiente en términos de tiempo para búsquedas repetidas.

En la práctica, el uso de estructuras de datos como mapas o índices invertidos es común para mejorar la eficiencia en la búsqueda de palabras en grandes conjuntos de documentos. La elección del algoritmo y la estructura de datos adecuados depende de los requisitos específicos de su aplicación y de la cantidad de documentos que necesita procesar.