

## Análisis de Algoritmos

**Actividad 6**

Leonardo Flores Torres

17 de diciembre de 2022

Programar el algoritmo de Horspool para la búsqueda de palabras clave para los siguientes casos:

- Caso 1:
  - Entrada: Palabra de búsqueda, y archivo de texto `.txt`
  - Salida: No se encuentra la palabra, o bien sí se encuentra y la posición donde comienza la palabra (renglón y columna).
- Caso 2:
  - Entrada: Palabra de búsqueda y múltiples archivos de texto `.txt` (en una carpeta).
  - Salida: La lista de los archivos ordenados por relevancia (de mayor a menor). Esto es, cuántas veces se encuentra la palabra en cada archivo, y si dos archivos tienen el mismo conteo entonces será más relevante aquel en el la razón de las ocurrencias de la palabra de búsqueda respecto al total de palabras del documento sea mayor.

**Solución:**

El algoritmo de Horspool, también conocido como el algoritmo Boyer-Moore-Horspool, fue hecho para la búsqueda de palabras dentro de texto. Esto puede verse como la búsqueda de una subcadena de texto dentro de una cadena de entrada. Ejemplos del uso de este algoritmo se muestran en el curso *Data structures and Algorithms* [1] y en las notas de clase basadas en el libro *Análisis de algoritmos* [2].

Este algoritmo busca en una cadena un patrón deseado, se va recorriendo la cadena hasta que se encuentra la primer letra del patrón y se guarda su posición en una tabla de ocurrencias, si el siguiente carácter de la cadena también coincide con el siguiente de la palabra buscada se vuelve a guardar su posición dentro de la tabla, y esto se repite hasta que todos los caracteres de la palabra clave son encontrados. Pero se deben realizar un par de consideraciones.

Para realizar esta búsqueda se sobreponen los caracteres  $s_i$  de la cadena con aquellos del patrón  $c_j$ ,

$$\begin{array}{ccccccc} s_0 & s_1 & s_2 & \dots & s_n \\ c_0 & c_1 & c_2 & & & & \end{array},$$

si hay un carácter de la cadena  $s_i$  que no esté dentro del patrón se mueve el patrón un número de espacios igual a su longitud y se realiza nuevamente la búsqueda, por ejemplo si  $s_0 = c_0$ ,  $s_2 = c_2$

pero  $s_1 \neq c_1$ ,

$$\begin{array}{cccccccc} s_0 & s_1 & s_2 & s_3 & s_4 & s_5 & \dots & s_n \\ & & & c_0 & c_1 & c_2 & & \end{array} .$$

El segundo caso es si un carácter  $s_i$  de la cadena se encuentra en el patrón pero no coincide con la posición actual en el patrón, entonces se mueve la cadena para hacerlos coincidir

$$\begin{array}{cccccccc} s_0 & s_1 & s_2 & s_3 & s_4 & s_5 & \dots & s_n \\ & & c_0 & c_1 & c_2 & & & \end{array} ,$$

por ejemplo, si  $s_4 \neq c_2$  pero  $s_4 = c_0$ ,

$$\begin{array}{cccccccc} s_0 & s_1 & s_2 & s_3 & s_4 & s_5 & s_6 & \dots & s_n \\ & & & c_0 & c_1 & c_2 & \dots & & \end{array} .$$

El tercer caso es similar al primero, el último carácter del patrón coincide con su carácter en la cadena pero no hay coincidencias con el resto de caracteres del patrón,

$$\begin{array}{ccccccc} s_0 & s_1 & s_2 & \dots & s_n \\ c_0 & c_1 & c_2 & & \end{array} ,$$

esto es, si  $s_2 = c_2$  pero  $s_0 \neq c_0$  y  $s_1 \neq c_1$ , entonces la cadena se desplaza por la longitud total del patrón como en el primer caso,

$$\begin{array}{cccccccc} s_0 & s_1 & s_2 & s_3 & s_4 & s_5 & \dots & s_n \\ & & c_0 & c_1 & c_2 & & & \end{array} .$$

El último caso es cuando el último carácter del patrón coincide con su respectivo carácter en la cadena pero hay otras coincidencias en el resto de caracteres

$$\begin{array}{ccccccc} s_0 & s_1 & s_2 & s_3 & s_4 & \dots & s_n \\ c_0 & c_1 & c_2 & c_3 & & & \end{array} ,$$

por ejemplo, si  $s_4 = c_3$  pero  $s_3 \neq c_2$  y  $s_4 = c_1$ ,  $s_3 = c_0$ ,

$$\begin{array}{cccccccc} s_0 & s_1 & s_2 & s_3 & s_4 & s_5 & s_6 & \dots & s_n \\ & & c_0 & c_1 & c_2 & c_3 & & & \end{array} ,$$

entonces el desplazamiento de la cadena es similar al segundo caso.

Primero se carga el módulo `MiniFinder` en `julia` usando un alias para llamar más rápidamente a las funciones contenidas dentro de este. La búsqueda se hace mediante la función `horspool` que detecta solamente la primera ocurrencia del patrón en una cadena como se muestra a continuación:

```
1 julia> using Revise; using MiniFinder; mf = MiniFinder;
2 julia> text = "El curso de analisis de algoritmos es interesante. Leonardo estudia para su
  ↳ curso."
3 "El curso de analisis de algoritmos es interesante. Leonardo estudia para su curso."
4 julia> pattern = "leo"
5 "leo"
```

```

6 julia> mf.horspool(lowercase(text), pattern)
7 51

8 julia> pattern = "curso"
9 "curso"

10 julia> mf.horspool(lowercase(text), pattern)
11 3

```

En el caso en que el patrón es igual a "curso", el algoritmo de Horspool no detecta la segunda ocurrencia. Finalmente lo que se desea para esta actividad es justamente esto, el poder contar cuantas ocurrencias hay del patrón en un texto, aunque para ello primero hay que diseñar una función que nos permita hacer esto en una línea de texto ya que el leer un archivo de texto .txt normalmente se hace por líneas. Esto se logró hacer mediante la función `occurrencesInLine` la cual recibe dos argumentos, el primero siendo una cadena de texto que representa una línea leída de alguna fuente y el segundo argumento el patrón de búsqueda.

Tomando el mismo texto guardado en la variable `text` definida anteriormente y realizando la búsqueda de los mismos patrones se obtiene lo siguiente:

```

12 julia> pattern = "leo"
13 "leo"

14 julia> occurrences = mf.occurrencesInLine(lowercase(text), pattern)
15 1-element Vector{Any}:
16 52

17 julia> text[occurrences[1]:end]
18 "Leonardo estudia para su curso."

19 julia> pattern = "curso"
20 "curso"

21 julia> occurrences = mf.occurrencesInLine(lowercase(text), pattern)
22 2-element Vector{Any}:
23 4
24 77

25 julia> text[occurrences[1]:end]
26 "curso de analisis de algoritmos es interesante. Leonardo estudia para su curso."

27 julia> text[occurrences[2]:end]
28 "curso."

```

El objetivo es poder observar como se realiza la detección de las ocurrencias. Por ejemplo, si se tiene una cadena de longitud  $n_c$ , siendo  $s_i$  los caracteres que conforman esa cadena, y se busca un patrón de longitud  $n_p$ , asumiendo que se encuentra la primera ocurrencia del patrón en la posición  $s_k$  donde  $0 \leq k \leq n_c - n_p$ , esto no asegura no haya otra ocurrencia del patrón en el resto de caracteres desde  $s_{k+1}$  hasta  $s_{n_c - n_p}$ . Se debe escribir una función que tome en cuenta lo ya mencionado y que aplique nuevamente la búsqueda del patrón en el resto de la cadena manteniendo al mismo tiempo

la posición relativa de la nueva cadena respecto a la cadena original.

cadena original	$s_0$	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$\dots$	$s_n$
cadena nueva				$s'_0$	$s'_1$	$s'_2$	$\dots$	$s'_n$
índices	0	1	2	3	4	5	$\dots$	$n$

La función que se encarga de esto es `occurrencesInLines`, como ya se había mostrado anteriormente. De manera similar, al ser ahora posible leer una línea se puede extender este concepto a leer un grupo de líneas, lo que es más, un grupo de líneas es un archivo de texto `.txt`. La función `occurrencesInFile` lee un archivo en este formato en un directorio indicado. Lo hace aplicando la función `occurrencesInLine` y guarda el índice de la línea en la que estaba y su ocurrencia en la línea actual en forma de tuplas  $(line_i, position_j)$ . Estas tuplas se van juntando en una lista donde se colectan todas las ocurrencias detectadas, y en el caso en que no se encuentre ninguna ocurrencia del patrón en el archivo indicado, la función imprime un mensaje de no haber encontrado ninguna y regresa una lista vacía. Ejemplo de esto se puede ver con un archivo de prueba, como se muestra a continuación:

```

29 julia> occurrences = mf.occurrencesInFile("../files/notabook_paragraphs_03.txt", "sed");
30 julia> occurrences
31 10-element Vector{Any}:
32 (1, 58)
33 (1, 168)
34 (1, 243)
35 (1, 383)
36 (3, 169)
37 (3, 243)
38 (3, 338)
39 (3, 952)
40 (3, 999)
41 (5, 307)

42 julia> occurrences = mf.occurrencesInFile("../files/notabook_paragraphs_03.txt", "desk");
43 No occurrences of "desk" where detected in "../files/notabook_paragraphs_03.txt".

44 julia> occurrences
45 Any[]

```

Los archivos de prueba para esta actividad se obtuvieron de forma aleatoria del sitio Lorem Ipsum<sup>1</sup> en el que se pueden generar textos con longitud arbitraria. Se decidió buscar los textos de esta manera para evitar complicar la limpieza de textos `.pdf` convertidos a `.txt`, además de que el propósito de esta actividad es solamente implementar la búsqueda.

Para completar el último punto de la actividad se tienen que contar de alguna manera las palabras que existen en un archivo y así tener la densidad de palabras. Las funciones mencionadas hasta ahora están todas relacionadas a la detección de las ocurrencias de una sola palabra usada como patrón dentro de una cadena o conjunto de cadenas, no para contar cuáles son las palabras en un archivo ni su frecuencia. Se buscó un acercamiento mas *naive* para resolver el problema del conteo.

<sup>1</sup><https://loremipsum.io/>

Para contar palabras se lee un archivo por líneas, cada línea es a su vez dividida por palabras a partir de la detección de ciertos símbolos de puntuación comunes en inglés como lo son " .,:;- ' ! ? " , además de que todas las letras en mayúsculas son convertidas a minúsculas. Al momento de leer una palabra ésta se agrega a un diccionario que mantiene la cuenta de esa palabra mientras detecta una ocurrencia de ella en el texto, y este proceso se repite hasta haber pasado por todas las palabras dentro del archivo. La función dedicada para esta tarea es `wordCounterInFile` la cual recibe como argumento el `path` para un archivo. Tomando el mismo archivo anteriormente usado como ejemplo se muestra el siguiente conteo:

```

46 julia> wordcount, totalcount = mf.wordCounterInFile("../files/notabook_paragraphs_03.txt");
47
48 julia> wordcount
49 Dict{String, Int64} with 129 entries:
50 "maecenas"    => 2
51 "quam"       => 4
52 "odio"       => 3
53 "ullamcorper" => 4
54 "tempus"     => 3
55 "et"         => 3
56 "faucibus"   => 2
57 "dictumst"   => 1
58 "rutrum"     => 1
59 "vitae"      => 4
60 "tempor"     => 1
61 :           => :
62
63 julia> totalcount
64 382

```

Lo que se muestra en `wordcount` son las palabras con sus respectivos conteos de ocurrencias, mientras que `totalcount` es el número total de palabras dentro del documento. Además, es posible saber cuál es el número de ocurrencias para una sola palabra, por ejemplo, para la misma anteriormente usada `"sed"`:

```

63 julia> get(words, "sed", -1)
64 10
65
66 julia> get(words, "desk", -1)
67 -1

```

La función `get` busca en el diccionario por llave (en este caso la palabra), y regresa el valor asociado a esa llave. En el caso en que no se encuentra esa llave entonces regresa un valor por defecto, en este caso es `-1`.

El procedimiento anteriormente descrito es fácilmente aplicable a un grupo de archivos en un directorio. Se escribió una función que identifique los archivos presentes en un directorio, aplique la función `mf.wordCounterInFile`, y que calcule el conteo del patrón y la densidad del patrón respecto al total de palabras para una futura comparación. Esta información se guarda en una lista de tuplas

ordenadas como (patterncount, density, filename) pero aquel documento en que el patrón no se haya detectado no se incluye en la lista. Ejemplo de esto se puede observar a continuación:

```

67 julia> collection = mf.wordCountInDirectory("../files/", "sed")
68 8-element Vector{Any}:
69 (10, 0.02617801047120419, "notabook_paragraphs_03")
70 (36, 0.029605263157894735, "notabook_paragraphs_10")
71 (96, 0.027126306866346424, "notabook_paragraphs_27")
72 (81, 0.02188006482982172, "notabook_paragraphs_29")
73 (87, 0.020928554245850373, "notabook_paragraphs_31")
74 (94, 0.02087960906263883, "notabook_paragraphs_34")
75 (101, 0.018703703703703705, "notabook_paragraphs_40")
76 (115, 0.01984126984126984, "notabook_paragraphs_45")

77 julia> collection = mf.wordCountInDirectory("../files/", "mollis")
78 7-element Vector{Any}:
79 (2, 0.001644736842105263, "notabook_paragraphs_10")
80 (3, 0.0008476970895733258, "notabook_paragraphs_27")
81 (10, 0.002701242571582928, "notabook_paragraphs_29")
82 (10, 0.0024055809477988932, "notabook_paragraphs_31")
83 (19, 0.004220346512661039, "notabook_paragraphs_34")
84 (9, 0.0016666666666666668, "notabook_paragraphs_40")
85 (6, 0.0010351966873706005, "notabook_paragraphs_45")

86 julia> collection = mf.wordCountInDirectory("../files/", "desk")
87 Any[]

```

Si la palabra buscada no se encuentra, el archivo no se lista, y si ningún archivo incluye el patrón en su texto, entonces la lista resultante es vacía. Ya que se tiene la información pertinente a cada documento es necesario ordenarlos de acuerdo al criterio mencionado en la actividad, este es, en orden descendente de acuerdo a su conteo de palabras. Para este propósito hay dos funciones que se pueden elegir de acuerdo a lo deseado, si solamente se quiere ordenar de acuerdo al número de ocurrencias o a la densidad del patrón respecto al total se puede utilizar la función `sortCollection`. Un ejemplo de esto se muestra a continuación al buscar el patrón "sed" en la colección de archivos usada hasta ahora:

```

88 julia> collection = mf.wordCountInDirectory("../files/", "sed")
89 8-element Vector{Any}:
90 (10, 0.02617801047120419, "notabook_paragraphs_03")
91 (36, 0.029605263157894735, "notabook_paragraphs_10")
92 (96, 0.027126306866346424, "notabook_paragraphs_27")
93 (81, 0.02188006482982172, "notabook_paragraphs_29")
94 (87, 0.020928554245850373, "notabook_paragraphs_31")
95 (94, 0.02087960906263883, "notabook_paragraphs_34")
96 (101, 0.018703703703703705, "notabook_paragraphs_40")
97 (115, 0.01984126984126984, "notabook_paragraphs_45")

98 julia> mf.sortCollection(collection, "wordcount")
99 8-element Vector{Any}:
100 (115, 0.01984126984126984, "notabook_paragraphs_45")
101 (101, 0.018703703703703705, "notabook_paragraphs_40")
102 (96, 0.027126306866346424, "notabook_paragraphs_27")

```

```

103 (94, 0.02087960906263883, "notabook_paragraphs_34")
104 (87, 0.020928554245850373, "notabook_paragraphs_31")
105 (81, 0.02188006482982172, "notabook_paragraphs_29")
106 (36, 0.029605263157894735, "notabook_paragraphs_10")
107 (10, 0.02617801047120419, "notabook_paragraphs_03")

108 julia> mf.sortCollection(collection, "density")
109 8-element Vector{Any}:
110 (36, 0.029605263157894735, "notabook_paragraphs_10")
111 (96, 0.027126306866346424, "notabook_paragraphs_27")
112 (10, 0.02617801047120419, "notabook_paragraphs_03")
113 (81, 0.02188006482982172, "notabook_paragraphs_29")
114 (87, 0.020928554245850373, "notabook_paragraphs_31")
115 (94, 0.02087960906263883, "notabook_paragraphs_34")
116 (115, 0.01984126984126984, "notabook_paragraphs_45")
117 (101, 0.018703703703703705, "notabook_paragraphs_40")

```

Aunque lo anterior es necesario para tener una idea de la importancia de nuestros archivos respecto a la ocurrencia del patrón no es suficiente para cumplir con el último inciso de la asignatura ¿Cómo diseñar un algoritmo de ordenamiento para ordenar de acuerdo a dos parámetros? Con esto en mente se buscó implementar un algoritmo sencillo [3], y además efectivo, que itera dos veces sobre la lista original intercambiando valores si detecta que uno es menor que otro, se adaptó para ajustarse a la lista de tuplas obtenida donde se guarda la información de los archivos y se agregó una condición extra para manejar el caso cuando el conteo de las palabras en un archivo es igual a otro:

```

1 function specialSort(collection)
2     # Remove from collection elements without the pattern, i.e., word count and
3     # density equal to -1
4     collection = filter(x → (x[1] != -1 && x[2] != -2), collection)
5     n = length(collection)
6
7     wordcount = map(x → x[1], collection)
8     density = map(x → x[2], collection)
9
10    for i in 1:n
11        for j in 1:n
12            if wordcount[i] < wordcount[j]
13                wordcount[i], wordcount[j] = wordcount[j], wordcount[i]
14                density[i], density[j] = density[j], density[i]
15                collection[i], collection[j] = collection[j], collection[i]
16            end
17            if wordcount[i] == wordcount[j]
18                newi, newj = density[i] < density[j] ? (j, i) : (i, j) # index by density
19                ↪ comparisson
20
21                wordcount[i], wordcount[j] = wordcount[newj], wordcount[newi]
22                density[i], density[j] = density[newj], density[newi]
23                collection[i], collection[j] = collection[newj], collection[newi]
24            end
25        end
26    end

```

Aunque si se aplica a la colección guardada en la variable `collection` del extracto de código inmediatamente arriba no habrá ningún cambio aparente

```
118 julia> mf.specialSort(collection)
119 8-element Vector{Any}:
120 (10, 0.02617801047120419, "notabook_paragraphs_03")
121 (36, 0.029605263157894735, "notabook_paragraphs_10")
122 (81, 0.02188006482982172, "notabook_paragraphs_29")
123 (87, 0.020928554245850373, "notabook_paragraphs_31")
124 (94, 0.02087960906263883, "notabook_paragraphs_34")
125 (96, 0.027126306866346424, "notabook_paragraphs_27")
126 (101, 0.018703703703703705, "notabook_paragraphs_40")
127 (115, 0.01984126984126984, "notabook_paragraphs_45")
```

Se considerará un caso en el que la primera entrada y la última de `collection` tienen el mismo conteo de palabras que la entrada para el archivo `"notabook_paragraphs_27"`,

```
118 julia> collection[1] = (96, 0.02617801047120419, "notabook_paragraphs_03")
119 (96, 0.02617801047120419, "notabook_paragraphs_03")

120 julia> collection[8] = (96, 0.01984126984126984, "notabook_paragraphs_45")
121 (96, 0.01984126984126984, "notabook_paragraphs_45")

122 julia> collection
123 8-element Vector{Any}:
124 (96, 0.02617801047120419, "notabook_paragraphs_03")
125 (36, 0.029605263157894735, "notabook_paragraphs_10")
126 (96, 0.027126306866346424, "notabook_paragraphs_27")
127 (81, 0.02188006482982172, "notabook_paragraphs_29")
128 (87, 0.020928554245850373, "notabook_paragraphs_31")
129 (94, 0.02087960906263883, "notabook_paragraphs_34")
130 (101, 0.018703703703703705, "notabook_paragraphs_40")
131 (96, 0.01984126984126984, "notabook_paragraphs_45")
```

Ahora, si se aplica la función de ordenamiento `specialSort` se esperaría que estas entradas con el mismo conteo de palabras sean ordenadas de acuerdo a su densidad justo como se muestra

```
118 julia> mf.specialSort(collection)
119 8-element Vector{Any}:
120 (36, 0.029605263157894735, "notabook_paragraphs_10")
121 (81, 0.02188006482982172, "notabook_paragraphs_29")
122 (87, 0.020928554245850373, "notabook_paragraphs_31")
123 (94, 0.02087960906263883, "notabook_paragraphs_34")
124 (96, 0.027126306866346424, "notabook_paragraphs_27")
125 (96, 0.02617801047120419, "notabook_paragraphs_03")
126 (96, 0.01984126984126984, "notabook_paragraphs_45")
127 (101, 0.018703703703703705, "notabook_paragraphs_40")
```



De esta manera termina la actividad de búsqueda y ordenamiento de textos de acuerdo al criterio de ocurrencias de un patrón en ellos, y al segundo criterio de ordenamiento sí es que algunos de ellos presentan el mismo conteo de ocurrencias.

## Apéndice

```

1  module MiniFinder
2
3  using OffsetArrays    # To create arrays with first index equal to zero.
4
5  #=
6  Important notes
7
8  Link to Horspool algorithm explanation
9  - http://www.cs.emory.edu/~cheung/Courses/253/Syllabus/Text/Matching-Boyer-Moore2.html
10
11  Link of available ascii characters:
12  - https://www.rapidtables.com/code/text/ascii-table.html
13
14  How to 0-base index in Julia?
15  - https://medium.com/analytics-vidhya/0-based-indexing-a-julia-how-to-43578c780c37
16
17  Normalize entry lines
18  julia> Base.Unicode.normalize("día? ~", stripmark=true, stripcc=true)
19
20  How to not use OffsetVectors?
21  How to use a dictionary for the occurrence table?
22  =#
23
24  function occurrenceTable(textstring)
25      # T → textstring
26
27      # textstring = isa(textstring, String) ? split(textstring, "") : textstring
28      textstring = split(textstring, "")
29      ascii = 256 # 128 # 256    # How many allowed ascii characters? There are 255 in total.
30
31      T = OffsetVector(textstring, 0:(length(textstring)-1))
32      lastocc = OffsetVector{Int, ASCII}(ones{Int, ASCII}(length(T), 256), 0:(length(T)-1)) # Initialize all to -1.
33
34      for i in 0:(length(T)-1)    # Antes solo estaba considerando un menos 1
35          lastocc[Int(T[i][1])] = i
36      end
37
38      return lastocc
39  end
40
41  function horspool(textstring, pattern)
42      # T → textstring
43      # P → pattern
44
45      # Dont even attempt it if the length of the textstring is smaller than
46      # the length of the pattern.
47      if length(textstring) < length(pattern)
48          return -1
49      end
50
51      lastocc = occurrenceTable(pattern)
52      textstring = split(textstring, "")

```

```

53     pattern = split(pattern, "")
54
55     T = OffsetVector(textstring, 0:length(textstring)-1)
56     P = OffsetVector(pattern, 0:length(pattern)-1)
57
58     n = length(T)
59     m = length(P)
60
61     i0 = 0
62
63     while i0 <= n - m
64         j = m - 1 # Start at the last char in the pattern
65
66         # When the last char in the pattern is found, iterate over all chars and
67         # compare to see if all match.
68         while Int(P[j][1]) == Int(T[i0 + j][1])
69             j = j - 1
70
71             if j < 0
72                 return i0 # Starts indexing at 0
73             end
74         end
75
76         i0 = i0 + (m - 1) - lastocc[Int(T[i0 + (m - 1)][1])]
77     end
78
79     return -1
80 end
81
82 function occurrencesInLine(line, pattern)
83     patternlength = length(pattern)
84     linelength = length(line)
85
86     line_ = line
87
88     occurrence = 0
89     currentpos = 1
90
91
92     while length(line_) >= patternlength
93         # while occurrence != -1
94         occurrence = horspool(line_, pattern)
95
96         if occurrence == -1
97             break
98         end
99
100        append!(occurrencelist, occurrence + currentpos)
101        currentpos = currentpos + occurrence + patternlength
102
103        if currentpos + patternlength > linelength
104            break
105        end
106

```

```

107     line_ = line[currentpos:end]
108     # println(line_)
109 end
110
111 return occurancelist
112 end
113
114 function occurrencesInFile(file, pattern)
115     occurlist = []
116
117     open(file, "r") do io
118         for (index, line) in enumerate(eachline(io))
119             line_ = lowercase(line)
120             occurrences = occurrencesInLine(lowercase(line_), pattern)
121
122             if !isempty(occurrences)
123                 # push!(occurlist, (index, occurrences))
124                 occurrencesintuples = map(x → (index, x), occurrences)
125                 append!(occurlist, occurrencesintuples)
126             end
127         end
128     end
129
130     if isempty(occurlist)
131         println("No occurrences of \"$(pattern)\" where detected in \"$(file)\".")
132     end
133     return occurlist
134 end
135
136 function occurrencesInDirectory(directory, pattern)
137     files = readdir(directory)
138     occurdict = Dict{String, Vector{Tuple{Int, Vector{String}}}}()
139
140     for file in files
141         filename = findlast('.', file) > x → file[begin:x-1] # remove extension from file name
142         occurrences = occurrencesInFile(directory * file, pattern)
143
144         occurdict[filename] = occurrences
145     end
146     return occurdict
147 end
148
149 function wordCountInFile(io)
150     countbyword = Dict{String, Int}{}
151     totalwordcount = 0
152     punctuationmarks = " .,:;-'!?"
153
154     for line in eachline(io), word in split(line, in(punctuationmarks)) # Add punctuation symbols as
155         ↪ needed
156         lword = lowercase(word)
157         countbyword[lword] = get(countbyword, lword, 0) + 1
158         totalwordcount += 1
159     end

```

```

160     countbyword, totalwordcount
161 end
162
163 function wordCountInFile_(io, counter=Dict{String,Int}())
164     for line in eachline(io), word in split(lowercase(line), !in('a':'z'))    # Add punctuation symbols as
165         ↪ needed
166         counter[word] = get(counter, word, 0) + 1
167     end
168     counter
169 end
170
171 function wordCountInDirectory(directory, pattern)
172     files = readdir(directory)
173     collection = []
174
175     for file in files
176         filename = findlast('.', file) ▷ x → file[begin:x-1]    # remove extension from file name
177         countbyword, totalwordcount = wordCountInFile(directory * file)
178         patterncount = get(countbyword, pattern, -1)
179         # density = !isequal(patterncount, -1) ? patterncount / totalwordcount : -1
180
181         if patterncount != -1
182             density = patterncount / totalwordcount
183             push!(collection, (patterncount, density, filename))
184         end
185     end
186
187     return collection
188 end
189
190 """
191     sortCollection(collection, sortby)
192
193     Sorts a vector of (wordcount, density, book) by wordcount or density in
194     descending order.
195 """
196 function sortCollection(collection, sortby)
197     options = Dict{"wordcount" ⇒ 1, "density" ⇒ 2}
198     option = get(options, sortby, -1)
199
200     if !isequal(option, -1)
201         return sort(collection, by = x → x[option], rev=true)
202     end
203 end
204
205 function specialSort(collection)
206     # Remove from collection elements without the pattern, i.e., word count and
207     # density equal to -1
208     collection = filter(x → (x[1] != -1 && x[2] != -2), collection)
209     n = length(collection)
210
211     wordcount = map(x → x[1], collection)
212     density = map(x → x[2], collection)

```

```

213
214     for i in 1:n
215         for j in 1:n
216             if wordcount[i] < wordcount[j]
217                 wordcount[i], wordcount[j] = wordcount[j], wordcount[i]
218                 density[i], density[j] = density[j], density[i]
219                 collection[i], collection[j] = collection[j], collection[i]
220             end
221             if wordcount[i] == wordcount[j]
222                 newi, newj = density[i] < density[j] ? (j, i) : (i, j) # index by density comparisson
223
224                 wordcount[i], wordcount[j] = wordcount[newj], wordcount[newi]
225                 density[i], density[j] = density[newj], density[newi]
226                 collection[i], collection[j] = collection[newj], collection[newi]
227             end
228         end
229     end
230
231     return collection
232 end
233
234 # function occurrenceDict(textstr)
235 #     lastocc = Dict()
236
237 #     for (i, char) in enumerate(textstr[begin:end-1])
238 #         lastocc[char] = i
239 #     end
240
241 #     return lastocc
242 # end
243
244 end # module MiniFinder

```

## Referencias

- [1] Shun Yan Cheung. The boyer-moore-horspool algorithm. <http://www.cs.emory.edu/~cheung/Courses/253/Syllabus/Text/Matching-Boyer-Moore2.html>, 2013. Visitado: 2022-11-25.
- [2] Homero V. Ríos Figueroa, Fernando M. Montes Gonzáles, Víctor R. Cruz Álvarez. Análisis de algoritmos. Universidad Veracruzana, 2013.
- [3] Stanley P. Y. Fung. Is this the simplest (and most surprising) sorting algorithm ever? *CoRR*, abs/2110.01111, 2021.
- [4] Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B. Shah. Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1):65–98, 9 2017.