

Sobre Matching en Hileras

Victor David Coto Solano

Diego Quirós Artiñano

Derek Rojas Mendoza

Universidad Nacional de Costa Rica

EIF-203: Estructuras Discretas Grupo 01-10am

Carlos Loria-Saenz

Ciclo I 2022



Resumen

String Matching son tipos de algoritmos esenciales para el mundo del siglo 21 debido a las aplicaciones variadas que este presenta. La importancia de poder entender y usar estos algoritmos va a ayudar en el ámbito laboral del futuro a corto y largo plazo. Nuestra investigación busca explicar de manera fácil 4 algoritmos diferentes para visualizar la eficiencia de este proceso.

Keywords: *String Matching, Fuerza Bruta, Knuth-Morris-Pratt, Boyer-Moore, Karp-Rabin, O-grande, ADN, Huntington*

Índice

| | |
|--------------------------------|-----------|
| Introducción General | 1 |
| <i>String Matching</i> | 1 |
| Algoritmos | 1 |
| Fuerza Bruta (Naive Algorithm) | 2 |
| Knuth-Morris-Pratt (KMP) | 4 |
| Algoritmo de Boyer-Moore | 7 |
| Algoritmo de Karp-Rabin | 11 |
| ADN | 16 |
| Referencias | 19 |

Índice de figuras

| | | |
|----|---|----|
| 1. | Fuerza Bruta (Apéndice fuerzaBruta.ipynb) | 3 |
| 2. | Knuth-Morris-Pratt (Apéndice kmp.ipynb) | 7 |
| 3. | Boyer-Moore (Apéndice boyermoore.ipynb) | 11 |
| 4. | Karp-Rabin (Apéndice rabinkarp.ipynb) | 15 |
| 5. | Todos los algoritmos comparados (Apéndice todos.ipynb) | 15 |
| 6. | Causas de Huntington (<i>Causas de la Enfermedad de Huntington</i> (2022)) | 17 |

Índice de tablas

| | | |
|----|---|----|
| 1. | Harry Potter 1 buscando Harry | 15 |
| 2. | ADN aleatorio buscando adeina, generada por <i>Random DNA Generator</i> (2022) (Apéndice pruebas.ipynb) | 16 |
| 3. | Secuencia de ejemplo sacada de <i>ENFERMEDAD DE HUNTINGTON. FENÓMENO DE ANTICIPACION GÉNICA: IMPORTANCIA DE LA CUANTIFICACIÓN DEL NÚ- MERO DE REPETCIONES CAG.</i> (s.f.), buscando Huntington (Apéndice pruebas.ipynb) | 17 |

Índice de algoritmos

| | | |
|----|---|----|
| 1. | Algoritmo de fuerza bruta | 2 |
| 2. | Algoritmo de Knuth-Morris-Pratt | 5 |
| 3. | Algoritmo de Boyer_Moore | 8 |
| 4. | Algoritmo de Karp-Rabin | 12 |

Introducción General

En esta investigación se cubrirá el tema de string matching y se explicarán, analizarán y construirán los 4 algoritmos solicitados por el SPEC. Se dará una breve introducción y se desarrollará cada algoritmo y se le realizará el análisis de tiempo visto en clase. Se hablará y se explicará brevemente sobre el ADN. Se hablará y se explicará brevemente la enfermedad de Huntington, esta misma fue seleccionada para realizar un ejemplo de la función de los algoritmos string matching y su importancia en la vida cotidiana, como por ejemplo; en este caso, evaluaremos analizando una cadena de ADN, si una persona padece dicha enfermedad.

String Matching

¿Qué es String Matching?

String Matching es cuando se agarra un patrón y se buscan todas las instancias de ese patrón dentro de un texto específico. String es el término que se usan para las cadenas de texto y matching la palabra en inglés que dice que dos cosas concuerdan, en este caso la hilera del patrón con la hilera del texto.

String Matching tiene varios usos en la vida real. Situaciones tan variadas como bases de datos musicales (Faro y Lecroq, 2013a) y a la misma vez se puede usar para detección de plagio, forense digital, chequeo de palabras al escribir, detección de intrusión y muchísimos otros casos posibles (*Applications of String Matching Algorithms - GeeksforGeeks* (2020)). Para esta investigación nos vamos a basar en secuencias de ADN como uso del String Matching.

Tipos de String Matching

String Matching se subdivide en dos categorías: exactas y semejantes. Las dos de las categorías tienen sus propios algoritmos diferentes. En esta investigación no solo le daremos enfoque a la categoría de algoritmos de String Matching exactos, sino que nos vamos a analizar 4 algoritmos de este, los cuales serán: Fuerza Bruta (el más básico), Knuth-Morris-Pratt, Boyer-Moore y finalmente el Rabin-Karp. A cada uno se le va a dar un enfoque mayor en la siguiente parte de la investigación, pero en general lo que los destaca es que los primeros tres son con revisión de caracteres y el último utiliza hashing (que convierte uno o varios elementos de entrada a una función en otro elemento).

Algoritmos

Introducción

En esta parte de la investigación se van a analizar los cuatro algoritmos previamente mencionados. Para esto se va a explicar en términos generales como sirve el programa, se va a generar un algoritmo y finalmente se va a analizar este algoritmo con los seis pasos vistos en clase.

Fuerza Bruta (Naive Algorithm)

Introducción al Algoritmo de Fuerza Bruta

El algoritmo de fuerza bruta es un algoritmo que analiza de izquierda a derecha revisando por caracteres si un patrón existe dentro de un texto. Se puede considerar el método más básico de String Matching porque en cada carácter revisa si el patrón se cumple.

Implementación del Algoritmo de Fuerza Bruta

(Con ayuda de *Programador CLICK* (s.f.), Sanz (2018) y Tello (s.f.))

Algoritmo 1 Algoritmo de fuerza bruta

```

1: procedure FUERZABRUTA((texto, patron))
2:    $n \leftarrow \text{len}(\text{texto})$ 
3:    $m \leftarrow \text{len}(\text{patron})$ 
4:   for  $i \leq (n - m)$  do                                ▷ Despues de  $n - m$  no puede ser el patron
5:     for  $j < m$  do                                       ▷ Evalua el patron caracter por caracter
6:       if  $\text{texto}[i + j] \neq \text{patron}[j]$  then           ▷ Evalua si los caracteres son los mismos
7:         break
8:       end if
9:     end for
10:    if  $j = m$  then                                       ▷ Si llega al final entonces existe y devuelve la posición
11:      return  $i$ 
12:    end if
13:  end for
14:  return -1                                              ▷ Si pasa por todo y no encuentra no existe devuelva -1
15: end procedure

```

Análisis del Algoritmo de Fuerza Bruta

Paso 1: Establecer el tamaño n de los datos

Hay dos variables que determinan el tamaño de los datos, estos serían n como la cantidad de caracteres del patrón y m la cantidad del texto. Entonces termina siendo
 $\text{len}(\text{patrón}) + \text{len}(\text{texto}) = n + m$

Paso 2: Determinar las operaciones de interés

Las operaciones de interés del algoritmo son las comparaciones de los caracteres. Estos tendrán valor de 1 cada uno.

Paso 3: Encontrar los casos base

Este algoritmo implica la combinación de las dos variables de tamaño un tiempo de recorrido del patrón y otro del texto.

$$T_{\text{patrón}}(0) = 0$$

$$T_{\text{patrón}}(1) = 1$$

$$T_{\text{patrón}}(n) = 1 + T_{\text{patrón}}(n - 1)$$

$$T_{\text{texto}}(0) = 0$$

$$T_{\text{texto}}(1) = 1$$

$$T_{\text{texto}}(m) = 1 + T_{\text{texto}}(m - 1)$$

y esto se junta para que de (porque para cada caracter del texto se evalúan todos los del patrón):

$$T_{\text{FuerzaBruta}}(n, m) = T_{\text{texto}}(m) * T_{\text{patrón}}(n)$$

Paso 4: Evaluando la ecuación recursiva

Como es una suma de unos

$$T_{\text{patrón}}(n) = 1 + 1 + 1 + 1 \dots + 1 + 1$$

$$T_{\text{patrón}}(n) = n$$

$$T_{\text{texto}}(m) = 1 + 1 + 1 + 1 \dots + 1 + 1$$

$$T_{\text{texto}}(m) = m$$

Con la ecuación que sacamos en el paso anterior y la evaluación de las sub-ecuaciones entonces llegamos a la conclusión de que:

$$T_{\text{FuerzaBruta}}(n, m) = m * n$$

Paso 5: O-grande

Con la ecuación de tiempo entonces concluimos que el O-grande del algoritmo es: $O(n * m)$

Paso 6

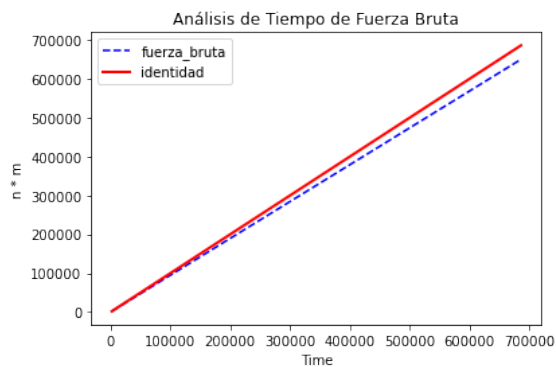


Figura 1

Fuerza Bruta (Apéndice fuerzaBruta.ipynb)

Knuth-Morris-Pratt (KMP)

Introducción al Algoritmo de Knuth-Morris-Pratt

Es un algoritmo que hace las comparaciones de izquierda a derecha, este realiza la búsqueda usando información basada en fallos previos obtenidos del patrón, esto se hace en una fase de “procesamiento” que crea una tabla de valores sobre su propio contenido. Esta tabla se crea para ver donde podría darse la siguiente coincidencia sin buscar más de 1 vez los caracteres del texto o cadena de caracteres donde se realiza la búsqueda, el algoritmo es de tiempo $O(n + m)$ siendo n = el tiempo de la fase de “searching” y m = el tiempo de la fase de procesamiento.

La tabla de fallos se encarga de evitar que cada carácter del texto sea analizado más de 1 vez, esto lo logra comparando el patrón consigo mismo para ver que partes se repiten. Este método guarda una lista con números le indican al algoritmo cuando debe devolverse desde la posición actual una vez que el patrón no coincide con el texto.

El texto y el patrón van avanzando simultáneamente mientras ambos coincidan, si una vez coinciden del todo pero la letra siguiente sigue cumpliendo con el patrón, entonces el algoritmo mueve el patrón 1 a la derecha, si no coinciden, entonces el patrón se empieza a devolver para intentar hacerlo coincidir con el texto.

Implementación del Algoritmo de Knuth-Morris-Pratt

(Con ayuda de *Knuth-Morris-Pratt algorithm* (2001), Colaboradores de los proyectos Wikimedia (2021) y *KMP Algorithm for Pattern Searching - GeeksforGeeks* (2022))

Algoritmo 2 Algoritmo de Knuth-Morris-Pratt

```

1: procedure KNUTH_MORRIS_PUTH((texto, patron))
2:    $n \leftarrow \text{len}(\text{texto})$ 
3:    $m \leftarrow \text{len}(\text{patron})$ 
4:    $\text{resultado} = \text{False}$ 
5:    $\text{listaDeIndices} = []$ 
6:    $\text{tablaDeFallo} = [0] * m$   $\triangleright$  Tabla que se va a usar para devoluciones en los procesamientos
7:    $i = 0$ 
8:    $j = 0$ 
9:   procedure PROCESAMIENTO((patron, m, tablaDeFallo))
10:     $\text{longitud} = 0$   $\triangleright$  longitud del sufijo del prefijo mas largo
11:     $i = 1$ 
12:    while  $i < m$  do
13:      if  $\text{patron}[i] == \text{patron}[\text{longitud}]$  then
14:         $\text{longitud}++ = 1$ 
15:         $\text{tablaDeFallo}[i] = \text{longitud}$ 
16:         $i++ = 1$ 
17:      else
18:        if  $\text{longitud}! = 0$  then
19:           $\text{longitud} = \text{tablaDeFallo}[\text{longitud} - 1]$ 
20:        else
21:           $\text{tablaDeFallo}[i] = 0$ 
22:           $i++ = 1$ 
23:        end if
24:      end if
25:    end while
26:  end procedure
27:  procedure BÚSQUEDA()
28:     $i = 0$ 
29:     $j = 0$ 
30:    while  $i < n$  do
31:      if  $\text{patron}[j] == \text{texto}[i]$ 
32:         $i++ = 1$ 
33:         $j++ = 1$ 
34:      end if
35:      if  $j == m$  then
36:         $\text{listaDeIndices}+ = [i - j]$ 
37:         $\text{resultado} = \text{True}$ 
38:         $j = \text{tablaDeFallo}[j - 1]$ 
39:      else if  $i < n \wedge \text{patron}[j]! = \text{text}[i]$  then

```

```

40:         if  $j! = 0$  then
41:              $j = \text{tablaDeFallo}[j - 1]$ 
42:         else
43:              $i+ = 1$ 
44:         end if
45:     end if
46: end while
47: end procedure
48: return listaDeIndices
49: end procedure

```

Análisis del Algoritmo de Knuth-Morris-Pratt

Paso 1: Establecer el tamaño n de los datos

El tamaño esta dado por la cantidad de elementos en el texto * la cantidad de elementos en el patron

$$n = \text{len}(\text{txt}) * \text{len}(\text{pat})$$

Paso 2: Determinar las operaciones de interés

Comparaciones entre Pat y Txt y construcción de lps[]

Suponemos el peor caso: todas son iguales a la mas grande de todas

Asumimos que la mas grande vale 1

Paso 3: Establecer la relacion de recurrencia para $T_{KMP}(n, m)$ con $n = \text{len}(\text{txt}) \wedge m = \text{len}(\text{Pat})$

Se obtuvo que:

$$T_{KMP}(n, m) = m \quad \text{sin} = 0$$

$$= (2 + t_{kmpS}(n - 1)) + (1 + t_{kmpP}(m - 1)) \quad \text{sin} > 0$$

con kmpS = KMP search y kmpP = MKP Processing

Paso 4: Resolver la ecuacion de $T_{KMP}(n, m)$ eliminando la recursion

$$T_{kmp}(n, m) = 2n + m$$

Paso 5: Determinar el orden de crecimiento asintotico de $T_{KMP}(n, m)$

$$T_{kmp}(n, m) \sim O(n + m)$$

con n = el tiempo de la fase de "searching" y m = el tiempo de la fase de procesamiento.

Paso 6

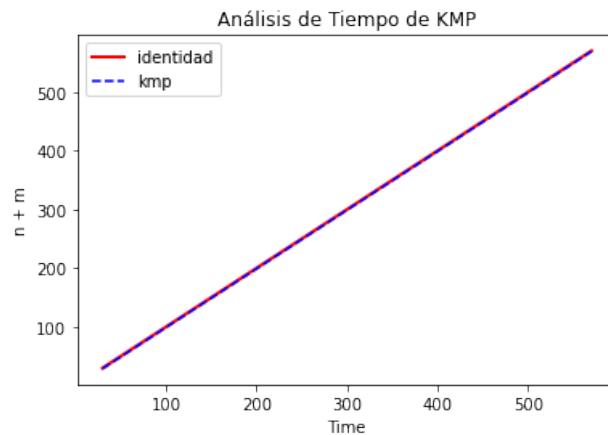


Figura 2

Knuth-Morris-Pratt (Apéndice kmp.ipynb)

Algoritmo de Boyer-Moore

Introducción del Algoritmo de Boyer-Moore

El algoritmo de Boyer-Moore realiza su búsqueda escaneando el patrón deseado desde la posición derecha hasta la posición izquierda. Este utiliza dos heurísticas para completar su cometido, llamadas 'Bad Character Rule' y 'Good Suffix Rule'. Para ambas heurísticas el algoritmo realiza un "preprocesamiento" del patrón, donde se elaboran dos tablas de interés:

- La tabla BC (Bad Character) basa su elaboración en el desplazamiento necesario para ir desde el carácter que se encuentra más a la derecha (del patrón) hasta la primera ocurrencia de los caracteres específicos que lo componen. Si
- La GS (Good Suffix) se elabora a partir de las coincidencias encontradas entre los sufijos del patrón y los caracteres restantes. Nos da el desplazamiento necesario, desde la izquierda, para encontrar dichas coincidencias. Estas pueden ser exactas (el prefijo aparece exactamente igual) o aproximadas (aparece una parte del sufijo).

Cuando el algoritmo va realizando las comparaciones entre el patrón y el texto y encuentre una discordancia, la decisión entre usar el desplazamiento recomendado por la tabla BC o la tabla GS dependerá netamente de cual de las dos le ofrezca un mayor "salto" o ventaja.

Implementación del Algoritmo de Boyer-Moore

(Con ayuda de *Boyer-Moore algorithm* (1997), Boyer y Moore (1977), Gusfield (1997), Faro y Lecroq (2013b) y *Exact Matching: Classical Comparison-Based Methods* (s.f.))

Algoritmo 3 Algoritmo de Boyer_Moore

```

1: procedure BOYER-MOORE((patron, texto))
2:    $sizeP = len(P)$ 
3:    $sizeT = len(T)$ 
4:    $boyerMooreBadChar = [0] * 256$       ▷ 256 es el número generalmente aceptado como
      alfabeto
5:   for  $0 \leq i < sizeP - 1$  do
6:      $boyerMooreBadChar[ord(patron[i])] = sizeP - i - 1$ 
7:   end for
8:    $suff = [0] * sizeP$ 
9:    $f = 0$ 
10:   $g = sizeP - 1$ 
11:   $suff[sizeP - 1] = sizeP$ 
12:  for  $sizeP - 2 \geq i > -1$  do
13:    if  $i > g \wedge suff[i + sizeP - 1 - f] < i - g$  then
14:       $suff[i] = suff[i + sizeP - 1 - f]$ 
15:    else
16:      if  $i < g$  then
17:         $g = i$ 
18:      end if
19:       $f = i$ 
20:      while  $g \geq 0 \wedge P[g] == P[g + sizeP - 1 - f]$  do
21:         $g = g - 1$ 
22:      end while
23:       $suff[i] = f - g$ 
24:    end if
25:  end for
26:   $boyerMooreGoodSuffix = [sizeP] * sizeP$ 
27:  for  $0 \leq i < sizeP$  do
28:    if  $suff[i] == i + 1$  then
29:      for  $0 \leq j < sizeP - 1 - i$  do
30:        if  $boyerMooreGoodSuffix[j] == sizeP$  then
31:           $boyerMooreGoodSuffix[j] = sizeP - 1 - i$ 
32:        end if
33:      end for
34:    end if
35:  end for
36:  for  $0 \leq i < sizeP - 1$  do
37:     $boyerMooreGoodSuffix[sizeP - 1 - suff[i]] = sizeP - 1 - i$ 
38:  end for
39:   $i = 0$ 

```

```

40:   $j = 0$ 
41:  while  $j \leq \text{size}T - \text{size}P$  do
42:     $i = \text{size}P - 1$ 
43:    while  $i \neq -1 \wedge \text{patron}[i] == \text{texto}[i + j]$  do
44:       $i = i - 1$ 
45:    end while
46:    if  $i < 0$  then
47:       $\text{print}(j)$ 
48:       $j += \text{boyerMooreGoodSuffix}[0]$ 
49:    else
50:       $j += \max(\text{boyerMooreGoodSuffix}[i], \text{boyerMooreBadChar}[\text{ord}(T[i + j]) -$ 
       $\text{size}P + 1 + i])$ 
51:    end if
52:  end while
53: end procedure

```

Análisis del Algoritmo de Boyer-Moore

Paso 1: Establecer el tamaño n de los datos

El tamaño de los datos está dado por $m + n$ (Siendo ' m ' el tamaño del texto y ' n ' el tamaño del patrón)

Paso 2: Determinar las operaciones de interés

Las operaciones van a variar dependiendo del módulo. Principalmente, utilizaremos de referencia aquellas relacionadas a comparaciones ('==', '!=') y asignaciones ('=')

Paso 3: Encontrar los casos base

Dividiremos esta sección en Preprocesamiento y Búsqueda, de la forma:

■ Preprocesamiento

- Heurística de caracteres malos (w es el tamaño del alfabeto)

$$T_{HCM}(0) = w \quad \text{si } n = 0$$

$$T_{HCM}(n) = 1 + T_{HCM}(n - 1) \quad \text{si } n > 0$$

- Heurística de sufijos buenos

$$T_{HSB}(0) = 0 \quad \text{si } n \leq 1$$

$$T_{HSB}(n) = 1 + (n - 1) + T_{HSB}(n - 1) \quad \text{si } n > 1$$

■ Búsqueda

$$T_{patrón}(0) = 0$$

$$T_{patrón}(1) = 1$$

$$T_{patrón}(n) = 1 + T_{patrón}(n - 1)$$

$$T_{texto}(0) = 0$$

$$T_{texto}(1) = 1$$

$$T_{texto}(m) = 1 + T_{texto}(m - 1)$$

$$T_{Búsqueda}(n, m) = T_{patrón}(n) * T_{texto}(m)$$

Eso entonces se junta para que la ecuación de tiempo del Boyer-Moore de:

$$T_{Boyer-Moore}(m, n) = T_{HCM}(n) + T_{HSB}(n) + T_{Búsqueda}(n, m)$$

Paso 4: Evaluando la ecuación recursiva

■ Preprocesamiento

- Heurística de caracteres malos

$$T_{HCM}(n) = 1 + 1 + 1 + \dots + 1 + w$$

$$T_{HCM}(n) = n + w$$

- Heurística de sufijos buenos

$$T_{HSB}(n) = 1 + (n - 1) + 1 + (n - 2) + 1 + (n - 3) \dots + 0$$

$$T_{HSB}(n) = n + (n - 1) + (n - 2) + (n - 3) \dots + 0$$

$$T_{HSB}(n) = \frac{n * (n + 1)}{2} = \frac{n^2 + n}{2}$$

■ Búsqueda

$$T_{patrón}(n) = 1 + 1 + 1 + 1 \dots + 1 + 1$$

$$T_{patrón}(n) = n$$

$$T_{texto}(m) = 1 + 1 + 1 + 1 \dots + 1 + 1$$

$$T_{texto}(m) = m$$

$$T_{Búsqueda}(n, m) = n * m$$

$$T_{Boyer-Moore}(n, m) = (n + w) + \left(\frac{n^2 + n}{2}\right) + n * m$$

Paso 5: O-grande

Dado a los pasos anteriores podemos concluir que el O-grande de Boyer-Moore es $O(w) + O(n^2) + O(n * m) \sim O(n * m)$. Es importante destacar que aunque en el peor de los casos es $n * m$, pero que es considerado el más eficiente de los 4 porque el O-grande normal es $O(\frac{n}{m})$ como lo describen Crochemore y Perrin (1991)

Paso 6

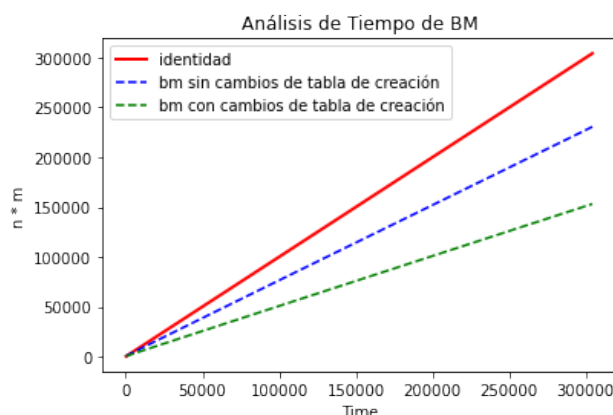


Figura 3

Boyer-Moore (Apéndice boyermoore.ipynb)

Este generamos dos porque en proceso de hacer pruebas nos percatamos que la generación de las tablas al ser uno o la otra no tenían que contarse todo el tiempo y creamos un if para que generara la más rápida.

Algoritmo de Karp-Rabin (o Rabin-Karp)

Introducción del Algoritmo de Karp-Rabin

El algoritmo de Karp-Rabin es el único de los que estamos analizando que usa el método de hashing. Para esto utiliza un número primo alto y una formula para sacar el hash value del patrón y de las secciones del texto que se están evaluando. Al encontrar un valor de hash que coincida entonces va a evaluar si la sección que está evaluando en el momento es igual que el patrón, carácter por carácter. Esta es la razón por la cual se utiliza un primo grande, porque al depender de residuos (primo todos los números van a salir modulo si mismo porque el primo solo es modulo 0 con si mismo) si el número evaluando es más grande que el primo entonces puede dar el mismo valor para lo que no debería.

Implementación del Algoritmo de Karp-Rabin

(Con ayuda de Cormen, Leiserson, Rivest, y Stein (2009), Hakak y cols. (2019), *Rabin-Karp Algorithm for Pattern Searching* - GeeksforGeeks (2022) y Karp-Rabin (s.f.))

Algoritmo 4 Algoritmo de Karp-Rabin

```

1: procedure KARP_RABIN((patron, texto))
2:    $n = \text{len}(\text{patron})$ 
3:    $m = \text{len}(\text{texto})$ 
4:    $d = 256$       ▷ Este es el alfabeto defecto que sale en varios analisis (caracteres alfabeto
                    inglés)
5:    $q = 33554393$       ▷ Cualquier número primo
                    sirve, pero preferiblemente alto porque los pequeños solo hacen que el algoritmo corra como
                    fuerza bruta porque más hashes concuerdan
6:    $h = d^{m-1} \bmod(q)$ 
7:    $\text{ValorHashPatron} = 0$ 
8:    $\text{ValorHashVentanaTexto} = 0$ 
9:    $\text{listaIndices} = []$ 
10:  for  $i = 0 < n$  do
11:     $\text{ValorHashPatron} = (d * \text{ValorHashPatron} + \text{patron}[i]) \bmod(q)$       ▷ Esta
                                operación sirve con un abecedario normal como tabla ascii, para usarlo en Python el accesor
                                se tiene que meter como parametro de ord()
12:     $\text{ValorHashVentanaTexto} = (d * \text{ValorHashVentanaTexto} + \text{texto}[i]) \bmod(q)$ 
13:  end for
     $j = 0$  ▷ definirla afuera para poder usarla dentro del for sin tener que redefinirla cada vez
    que empiece el for otra vez
14:  for  $i = 0 \leq m - n$  do
15:    if  $\text{ValorHashPatron} == \text{ValorHashVentanaTexto}$  then ▷ Solo hacer fuerza Bruta cuando
    los valores hash concuerdan
16:      for  $j = 0 < n$  do
17:        if  $\text{patron}[j] \neq \text{texto}[i + j]$  then      ▷ Si la fuerza bruta se incumple salga
18:          break
19:        else
20:           $j++ = 1$ 
21:        end if
22:      end for
23:      if  $j == n$  then

```

```

24:         listaIndices+ = [i]           ▷ Al llegar al final de la fuerza bruta registra el índice
25:     end if
26: end if
27: if i < m - n then
28:     ValorHashVentanaTexto = (d*(ValorHashVentanaTexto - texto[i]*h) + texto[i +
n])mod(q)
29:     if ValorHashVentanaTexto < 0 then   ▷ GeeksForGeeks recomienda en caso de
que den hashes negativos
30:         ValorHashVentanaTexto+ = q
31:     end if
32: end if
33: end for
34: return listaIndices
35: end procedure

```

Análisis del Algoritmo de Karp-Rabin

Paso 1: Establecer el tamaño n de los datos

Las dos variables que varían el tamaño de datos es la cantidad de caracteres del patrón n y la del texto m . Entonces termina siendo $\text{len}(\text{patrón}) + \text{len}(\text{texto}) = n + m$

Paso 2: Determinar las operaciones de interés

Las operaciones de interés son las comparaciones tanto del valor de hash como el carácter y van a contar como 1.

Paso 3: Encontrar los casos base

El algoritmo tiene tres sub-ecuaciones el de tiempo que tarda comparar los caracteres de patrón, los caracteres de texto y las comparaciones de los hashes.

$$T_{\text{patrón}}(0) = 0$$

$$T_{\text{patrón}}(1) = 1$$

$$T_{\text{patrón}}(n) = 1 + T_{\text{patrón}}(n - 1)$$

$$T_{\text{texto}}(0) = 0$$

$$T_{\text{texto}}(1) = 1$$

$$T_{\text{texto}}(m) = 1 + T_{\text{texto}}(m - 1)$$

$$T_{\text{hashes}}(0) = 0$$

$$T_{\text{hashes}}(1) = 1$$

$$T_{\text{hashes}}(m) = 1 + T_{\text{hashes}}(m - 1)$$

$$T_{Karp-Rabin}(n, m) = T_{hashes}(m) + T_{texto}(m) * T_{patrón}(n)$$

Se hace de esta manera porque siempre se evalúan los hashes y después por aparte se hace un fuerza bruta de la sección.

Paso 4: Evaluando la ecuación recursiva

$$T_{patrón}(n) = 1 + 1 + 1 + 1... + 1 + 1$$

$$T_{patrón}(n) = n$$

$$T_{texto}(m) = 1 + 1 + 1 + 1... + 1 + 1$$

$$T_{texto}(m) = m$$

$$T_{hashes}(m) = 1 + 1 + 1 + 1... + 1 + 1$$

$$T_{hashes}(m) = m$$

$$T_{Karp-Rabin}(n, m) = m + m * n$$

Paso 5: O-grande

Dado a lo que evaluamos en los últimos pasos se sabe que O-grande es

$$O(m) + O(m * n) \sim O(m * n)$$

Es importante notar que ese tiempo ocurre cuando todos los caracteres son lo mismo entonces tiene que evaluar todos los caracteres. El tiempo normal es $O(m) + O(cn) \sim O(m + n)$, c constante.

Paso 6

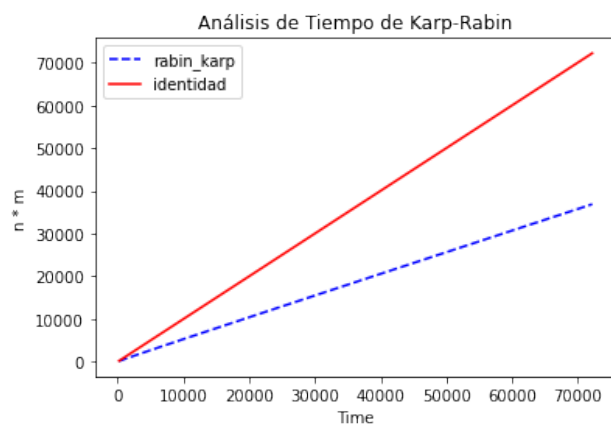


Figura 4

Karp-Rabin (Apéndice rabinkarp.ipynb)

Comparación de todos

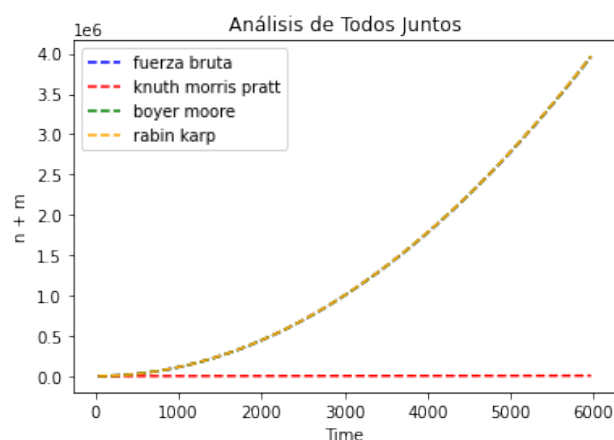


Figura 5

Todos los algoritmos comparados (Apéndice todos.ipynb)

Con esto podemos concluir que poniendo todos los algoritmos en el peor de los casos (todos los caracteres iguales) el kmp es el más eficiente.

Además decidimos probar los tres con un texto grande (buscar "Harry" dentro del primer libro de Harry Potter) y correrlo con el timeit de Python: (ver apéndice de pruebas.ipynb para más info)

| Fuerza Bruta | Knuth-Morris-Pratt | Boyer-Moore | Karp-Rabin |
|---------------------|---------------------|----------------------|----------------------|
| 92.39191730198218 s | 93.05609979297151 s | 138.52656350799953 s | 170.72219764100737 s |

Cuadro 1

Harry Potter 1 buscando Harry

ADN

Importancia del String Matching con el ADN

El Ácido desoxirribonucleico o ADN es un ácido nucleico compuesto de 4 nucleótidos, estos son la adenina (A), timina (T), guanina (G) y citosina (C), además de estos 4 nucleótidos, existe un 5 que es el uracilo (U), pero este se encuentra en el ARN reemplazando a la timina (T). Este contiene las instrucciones y la información genética usadas en el desarrollo y funcionamiento de todo ser vivo, incluyendo algunos virus, además, es el responsable de la herencia genética, como el color de ojos, el color de piel, el cabello y todo tipo de cosas que podemos heredar de nuestros padres o antepasados, sin embargo, también es responsable de heredar ciertas enfermedades genéticas, como por ejemplo, la fibrosis quística, la hemofilia, la enfermedad de Huntington, etc (*Las enfermedades hereditarias más comunes* | *Biolinks* (2022)).

Con esto se buscar una base nitrogenada específica, hicimos pruebas con los cuatro algoritmos con una secuencia de ADN aleatoria generada en *Random DNA Generator* (2022), y le corrimos con `timeit.repeat()` Para efectos de esta investigación, hemos decidido utilizar la

| | Fuerza Bruta | Knuth-Morris-Pratt | Boyer-Moore | Karp-Rabin |
|---|-----------------------|-----------------------|----------------------|-----------------------|
| H | 0.21230525700957514 s | 0.21744267400936224 s | 0.3150837090215646 s | 0.36325935399509035 s |

Cuadro 2

ADN aleatorio buscando adenina, generada por Random DNA Generator (2022) (Apéndice pruebas.ipynb)

enfermedad de Huntington como ejemplo de string matching en una cadena de ADN.

La enfermedad de Huntington es una enfermedad hereditaria que da la instrucción al cuerpo de producir una proteína llamada Huntingtina (HTT). Aunque se desconoce la función de dicha proteína, se cree que juega un papel importante en las neuronas. Esta enfermedad es provocada por una mutación en un segmento del ADN conocido como una repetición del trinucleótido CAG (citosina, adenina y guanina). Este segmento CAG normalmente se repite de 10 a 35 veces en personas sanas, sin embargo, en personas con la enfermedad de Huntington, dicho segmento se repite de 36 a 120 veces.

Como se puede ver en la siguiente imagen:

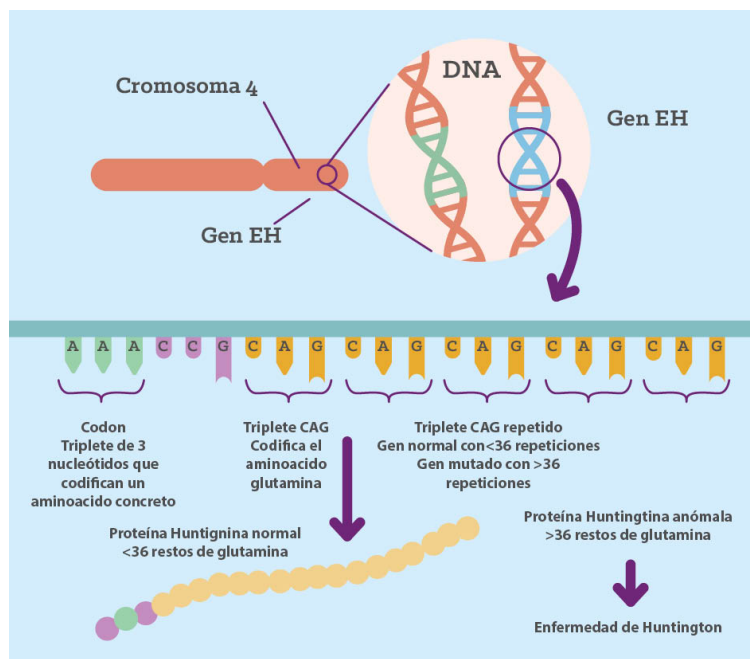


Figura 6

Causas de Huntington (Causas de la Enfermedad de Huntington (2022))

Para efectos del string matching, vamos a comprobar cuantas veces se repite el segmento CAG en la cadena de ADN para comprobar si una persona padece de esta enfermedad.

| | Fuerza Bruta | Knuth-Morris-Pratt | Boyer-Moore | Karp-Rabin |
|----------------|--------------|--------------------|-------------|------------|
| Sin Huntington | 0.0561 s | 0.0558 s | 0.089 s | 0.0815 s |
| Con Huntington | 0.0837 s | 0.0714 s | 0.1077 s | 0.1055 s |

Cuadro 3

Secuencia de ejemplo sacada de ENFERMEDAD DE HUNTINGTON. FENÓMENO DE ANTICIPACION GÉNICA: IMPORTANCIA DE LA CUANTIFICACIÓN DEL NÚMERO DE REPETICIONES CAG. (s.f.), buscando Huntington (Apéndice pruebas.ipynb)

Conclusión

De esta investigación aprendimos, investigamos, estudiamos e intentamos entender lo mejor posible sobre algunos de los algoritmos de string matching existentes. Se profundizó sobre definiciones y se explicaron los 4 algoritmos solicitados. Hicimos análisis de tiempo de los algoritmos y aprendimos acerca de la prueba experimental en el campo de la informática. Se habló sobre el ADN y la enfermedad de Huntington y se aplicó un ejemplo sobre dicha enfermedad para probar los algoritmos. Aprendimos sobre la importancia que pueden tener los algoritmos de string matching en la vida cotidiana, tanto en algo tan complejo como detectar enfermedades en cadenas de ADN como algo tan simple como buscar una palabra en un texto. En el proceso aprendimos más sobre diferentes conceptos relacionados con string matching, aprendimos más sobre Python y sobre los algoritmos solicitados. Creamos implementaciones de dichos algoritmos utilizando diferentes fuentes y pseudocódigos e intentamos darles distintos toques personales. Sin embargo, consideramos que aun tenemos un largo camino por recorrer

dado que algunos de los datos “teóricos” fueron muy diferentes a los prácticos, pero igualmente le dedicamos mucho esfuerzo..

Referencias

- Applications of String Matching Algorithms - GeeksforGeeks.* (2020, septiembre). Descargado de <https://www.geeksforgeeks.org/applications-of-string-matching-algorithms> ([Online; accessed 7. May 2022])
- Boyer, R. S., y Moore, J. S. (1977, oct). A fast string searching algorithm. *Commun. ACM*, 20(10), 762–772. Descargado de <https://doi.org/10.1145/359842.359859> doi: 10.1145/359842.359859
- Boyer-moore algorithm.* (1997). Descargado de <https://www-igm.univ-mlv.fr/~lecroq/string/node14.html>
- Causas de la Enfermedad de Huntington.* (2022, julio). Descargado de <https://rochepacientes.es/enfermedad-huntington/causas.html>
- Colaboradores de los proyectos Wikimedia. (2021, marzo). *Algoritmo Knuth-Morris-Pratt - Wikipedia, la enciclopedia libre.* Descargado de https://es.wikipedia.org/w/index.php?title=Algoritmo_Knuth-Morris-Pratt&oldid=133692421
- Cormen, T. H., Leiserson, C. E., Rivest, R. L., y Stein, C. (2009). *Introduction to algorithms* (3.^a ed.). London, England: MIT Press.
- Crochemore, M., y Perrin, D. (1991). *Two-way string-matching.* Descargado de <https://igm.univ-mlv.fr/~mac/Articles-PDF/CP-1991-jacm.pdf>
- Enfermedad de huntington. fenómeno de anticipacion gÉnica: Importancia de la cuantificaciÓn del nÚmero de repeticiones cag.* (s.f.). Descargado de <https://zaguan.unizar.es/record/47796/files/TAZ-TFG-2015-1112.pdf>
- Exact matching: Classical comparison-based methods.* (s.f.). Descargado de <http://pages.di.unipi.it/pisanti/DIDATTICA/patternmatching2.pdf>
- Faro, S., y Lecroq, T. (2013a). The exact online string matching problem: A review of the most recent results. *ACM computing surveys*, 45(2), 1–42.
- Faro, S., y Lecroq, T. (2013b, marzo). The exact online string matching problem: A review of the most recent results. *ACM Comput. Surv.*, 45(2), 1–42. doi: 10.1145/2431211.2431212
- Gusfield, D. (1997). Exact matching: Classical comparison-based methods. En *Algorithms on strings, trees, and sequences: Computer science and computational biology* (p. 16–34). Cambridge University Press. doi: 10.1017/CBO9780511574931.004
- Hakak, S. I., Kamsin, A., Shivakumara, P., Gilkar, G. A., Khan, W. Z., y Imran, M. (2019). Exact string matching algorithms: Survey, issues, and future research directions. *IEEE access*, 7, 69614–69637.
- Karp-Rabin, R. S. (s.f.). *String searching.* Descargado de <https://www.cs.princeton.edu/courses/archive/fall04/cos226/lectures/string.4up.pdf>
- KMP Algorithm for Pattern Searching - GeeksforGeeks.* (2022, junio). Descargado de <https://www.geeksforgeeks.org/kmp-algorithm-for-pattern-searching>
- Knuth-Morris-Pratt algorithm.* (2001, junio). Descargado de <http://www-igm.univ-mlv.fr/~lecroq/string/node8.html#SECTION0080>
- Las enfermedades hereditarias más comunes | Biolinks.* (2022, julio). Descargado de <https://biolinksperu.com/blog/enfermedades-hereditarias-comunes>
- Programador click.* (s.f.). Descargado de <https://programmerclick.com/article/21021656414/%7D>
- Rabin-Karp Algorithm for Pattern Searching - GeeksforGeeks.* (2022, junio). Descargado de <https://www.geeksforgeeks.org/rabin-karp-algorithm-for-pattern-searching>
- Random DNA Generator.* (2022, julio). Descargado de <http://www.faculty.ucr.edu/~mmaduro/random.htm>

- Sanz, B. C. (2018). *Algoritmos de búsqueda de subcadenas para encontrar semejanzas en cadenas numéricas* (Tesis Doctoral no publicada).
- Tello, E. A. R. (s.f.). *Algoritmos de fuerza bruta*. Descargado de <https://www.tamps.cinvestav.mx/~ertello/algorithms/sesion06.pdf>