

Implementación del algoritmo InexRecur en Python y comparación con implementación en R.*

Ignacio Amat Hernández

Ángela Llopis Hernández

Estrella Ballester Hoyo

3 de mayo de 2020

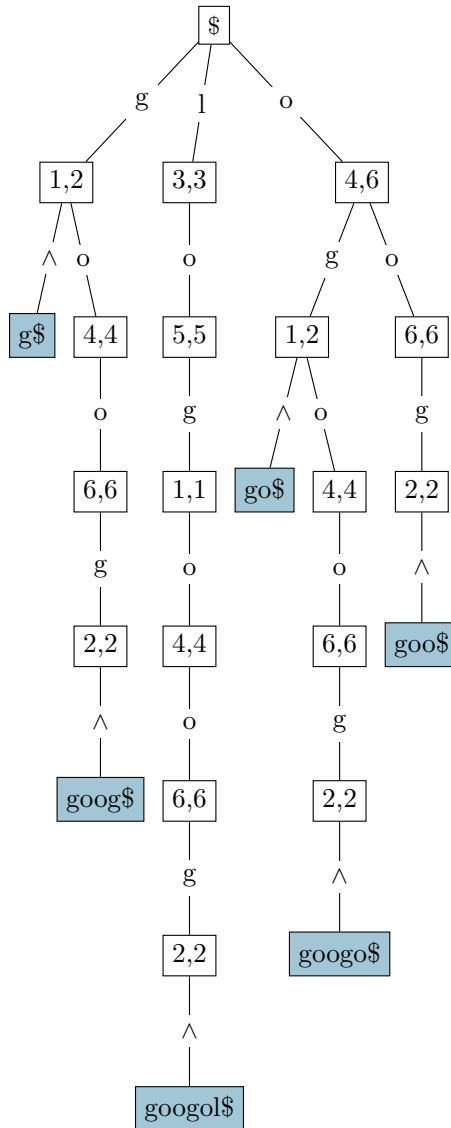


Figura 1: Árbol de prefijos con intervalos SA de la palabra “googol\$”.

*Grado en Ingeniería Biomédica, Escuela Técnica Superior de Ingenieros Industriales, Valencia, España.

Índice

1	Introducción	4
2	Métodos	5
2.1	Precálculo	5
2.1.1	Árbol de prefijos	5
2.1.2	Vector de sufijos $B[i]$ para la cadena de referencia X	5
2.1.3	Vector $C(a)$	5
2.1.4	Vector $O(a,i)$	6
2.1.5	Vector de sufijos $B'[i]$ para la cadena de referencia invertida X'	6
2.1.6	Vector $O'(a,i)$ para $B'[i]$	6
2.2	Búsqueda inexacta	6
2.2.1	Función <code>CalculateD</code>	6
2.2.2	Función de búsqueda inexacta recursiva <code>InexRecur</code>	7
2.2.3	Función <code>InexSearch</code>	8
3	Materiales	8
4	Resultados	9
5	Discusión	15
6	Conclusiones	15
7	Anejo	19

Índice de figuras

1	Árbol de prefijos con intervalos SA de la palabra “googol\$”	1
2	Ejemplo de la creación del vector de sufijos	5
3	Árbol de prefijos con intervalos SA de la palabra “aaaabbbbbb\$”	16
4	Árbol de prefijos con intervalos SA de la palabra “bananas\$”	17
5	Árbol de prefijos con intervalos SA de la palabra “cgatagtcgggatggttgccg\$”	18

Índice de cuadros

1	Obtención de los vectores $C(a)$ y $O(a,i)$ dado W y X	6
2	Algoritmo de la función <code>CALCULATED</code> . [2]	7
3	Algoritmo de la función <code>INEXRECUR</code> . [2]	7
4	Algoritmo de la función <code>INEXSEARCH</code> . [2]	8
5	Traza de <code>INEXRECUR</code> con $X = \text{“googol$”}$, $W = \text{“gol”}$, $z = 0$ en C y Python	9

6	Traza de INEXRECUR con $X = \text{"googol\$"}, W = \text{"gol"}, z = 0$ en R	9
7	Traza de INEXRECUR con $X = \text{"googol\$"}, W = \text{"goog"}, z = 0$ en C y Python	10
8	Traza de INEXRECUR con $X = \text{"googol\$"}, W = \text{"goog"}, z = 0$ en R	10
9	Traza de INEXRECUR con $X = \text{"googol\$"}, W = \text{"gool"}, z = 1$ en C y Python	11
10	Traza de INEXRECUR con $X = \text{"googol\$"}, W = \text{"gool"}, z = 1$ en R	12
11	Tiempos de "CPU".	13
12	Tiempo de "pared" según la utilidad <i>bench</i>	13
13	Memoria usada medido desde el script.	14
14	Memoria usada medido desde fuera del script.	14

Listings

1	Python <code>inexrecur_clean.py</code>	19
2	R <code>inexrecur_clean.R</code>	19
3	C <code>inexrecur_clean.c</code>	21
4	C <code>inexrecur_clean.c</code> (cont.)	21

Resumen

Debido a la aparición de nuevas técnicas de secuenciación masiva (NGS) en la bioinformática, existen actualmente numerosas herramientas para la secuenciación. Algunas de ellas han sido diseñadas para resolver el problema de alineamiento, entre ellas el Burrows-Wheeler Aligner (BWA), un alineador diseñado para alinear de secuencias a partir de un genoma de referencia [1]. BWA se basa en la transformada de Burrows-Wheeler (BWT) y tiene la capacidad de realizar búsquedas inexactas (permitiendo huecos y errores) de forma recursiva mediante el método InexRecur [2].

Hemos programado el algoritmo de recursión InexRecur en Python para el alineamiento de secuencias cortas con secuencias de referencia de gran tamaño y hemos realizado diversas pruebas comparando con la implementación del mismo algoritmo en R. Hemos obtenido las trazas en Python, R y C, los tiempos de CPU, la memoria requerida y los árboles de prefijos.

Los resultados obtenidos muestran que la implementación de la función de búsqueda recursiva inexacta en Python proporciona un menor tiempo de procesamiento y menor uso de memoria, además de poseer una sintaxis más intuitiva para el usuario en comparación con R.

1. Introducción

El campo de la bioinformática es una rama de la informática combinada con la biología, la química, la estadística y muchos otros ámbitos que trata de recopilar, almacenar, analizar, manipular, etc. datos de carácter biológico con el fin de comprender el funcionamiento de los organismos vivos, su origen y poder hallar así principios universales relativos a las ciencias de la vida [1]. Es un ámbito de gran interés para la industria científica pero que ha de manejar un gran número de datos (normalmente una máquina de secuenciación Illumina/Solexa puede producir de 50 a 200 millones de lecturas 32 a 100 pares de bases en una sola ejecución [2]) y con ello la bioinformática ha experimentado en los últimos años una gran evolución; cada vez secuenciar el genoma es más asequible y esto es gracias a las denominadas New Generation Sequencing (NGS) [1], nuevas técnicas más potentes y menos costosas computacionalmente. Dichas técnicas emplean algoritmos de alineamiento que permiten llevar a cabo una secuenciación masiva del ADN [2] con el fin de conocer sus secuencias de nucleótidos y poder compararlas entre individuos de las mismas o diferentes especies, así como hallar el origen evolutivo, la ortología u homología genética... Existen diversos relativos a la secuenciación genómica que las NGS tratan de abordar, entre ellos la resecuenciación, el mapeo o alineación y el ensamblado “de novo” del genoma[1].

En el presente trabajo nos centramos en el problema de la alineación que trata de, a partir de lecturas de un genoma desconocido, compararlas con un genoma de referencia conocido. Para dicha tarea existen numerosas técnicas que permiten resolver alineamientos pero no todas ellas funcionan de la misma forma. Entre los alineadores más empleados en la actualidad encontramos Bowtie y BWA. Nosotros nos centraremos en este último.

Burrows-Wheeler Aligner (BWA) es un alineador basado en la transformada de Burrows-Wheeler (BWT) que permite, a partir de un genoma de referencia, compararlo con secuencias cortas realizando búsquedas recursivas. Además, a diferencia de otros alineadores como Bowtie, BWA permite errores y gaps (huecos) en las secuencias, que son mutaciones genéticas muy comunes [1]. De esta forma BWA nos ofrece la posibilidad de realizar búsquedas exactas (sin errores ni gaps) o inexactas (con errores y gaps)[3]. En concreto, InexRecur es el algoritmo recursivo empleado por BWA para la búsqueda inexacta. Esta función es de gran utilidad a la hora de alinear secuencias de genoma muy largas, puesto que existe una gran probabilidad de haberse producido errores y gaps. En cambio Bowtie se suele emplear más a menudo en secuencias de menor tamaño en las que se busca una alta correspondencia.

Por otro lado, en cuanto a lenguajes de programación existen un gran número en el mercado. Y en concreto, para el análisis de datos Python y R son los más empleados para ello. En el presente artículo presentaremos el algoritmo en InexRecur en el lenguaje de programación Python y la compararemos con su implementación en R.

2. Métodos

2.1. Precálculo

2.1.1. Árbol de prefijos

El árbol de prefijos es una estructura con forma de árbol obtenida a partir de una cadena de texto de referencia X en la que cada carácter es referenciado por una arista. De esta manera, concatenando los caracteres de una arista de un nodo a la raíz se obtiene un substring único de X o cadena representada por el nodo. Cabe destacar que el árbol de prefijos de X es igual al árbol de sufijos de X invertida (X'). Esta herramienta facilita la búsqueda de, dado un substring W , saber si pertenece a X , pues solo habrá que buscar el nodo que representa a W . Podemos observar un ejemplo de un árbol de prefijos en la Fig. 1, 3, 4, 5.

2.1.2. Vector de sufijos ($B[i]$) para la cadena de referencia X

El primer paso de la búsqueda inexacta en Burrows-Wheeler es la creación del vector de sufijos $B[i]$ dada una cadena de texto de referencia X . $B[i]$ será una cadena de texto con los mismos caracteres que X pero ordenada lexicográficamente [4]. En la Figura 1 mostramos un ejemplo sencillo de la construcción de $B[i]$ a partir de un X de referencia. Hemos de recalcar que en la práctica, X sería una cadena de texto con un gran número de caracteres (nucleótidos) y W correspondería a una secuencia que queramos alinear con X .

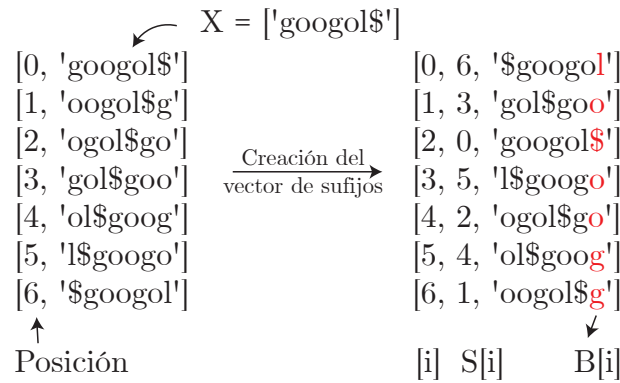


Figura 2: Ejemplo de la creación del vector de sufijos.

2.1.3. Vector $C(a)$

El vector $C(a)$ en Burrows-Wheeler se obtiene para observar si W es un substring perteneciente a X , y representa el número de caracteres en $X[0, n-2]$ lexicográficamente inferiores a a . En la Figura 2 mostramos un ejemplo del vector $C(a)$.

2.1.4. Vector $O(a,i)$

Otro vector implicado en BWA es $O(a,i)$ que indica el número de ocurrencias de un símbolo a en $B[0,i]$. Además, en InexRecur se puede realizar búsquedas inexactas de manera que, dado un W no perteneciente a X pero con un número de errores menores a un umbral z , podemos encontrar secuencias similares en el árbol de prefijos. En la Figura 2 podemos observar un ejemplo de $O(a,i)$.

X = 'googol\$'		B[i] = 'lo\$ooogg'	
W = 'lol'		C = (3, 2, 3)	
O(a,i)			
i	l	o	
0	1	0	
1	1	1	
2	1	1	
3	1	2	
4	1	3	
5	1	3	
6	1	3	

Cuadro 1: Obtención de los vectores $C(a)$ y $O(a,i)$ dado W y X .

2.1.5. Vector de sufijos $B'[i]$ para la cadena de referencia invertida X'

Siguiendo con el ejemplo de las Figuras 1 y 2, podríamos obtener a partir de $X = \text{'googol$'}$ su cadena invertida $X = \text{'$logooog'}$. Y el cálculo de $B'[i]$ de X' se realizará de forma análoga al de $B[i]$ para un string X . Como hemos mencionado anteriormente, el árbol de prefijos de X es equivalente al árbol de sufijos de X' .

2.1.6. Vector $O'(a,i)$ para $B'[i]$

Una vez obtenido $B'[i]$ calculamos su función correspondiente $O'(a,i)$ (Ejemplo en la Cuadro 1). Este es el último paso del precálculo, a partir del cual ya podremos llevar a cabo la búsqueda inexacta.

2.2. Búsqueda inexacta

2.2.1. Función CalculateD

Previamente a comenzar con la función InexRecur deberemos implementar la función CalculateD para obtener el vector D . Este vector hace referencia al número de fallos realizados hasta el momento, a variable z .

El parámetro de entrada a la función será el substring W , y el parámetro de salida será el propio vector D . La lógica llevada a cabo en la función será la mostrada en la Figura 3. El algoritmo establecerá unos valores para el umbral z y para el primer parámetro de W , j , e irá recorriendo en un bucle la subcadena W comprobando si pertenece o no a X . Finalmente asignará el valor de z a D y devolverá este último vector.

```

1  CALCULATED( $W$ )
2       $z \leftarrow 0$ 
3       $j \leftarrow 0$ 
4      for  $i = 0$  to  $|W| - 1$  do
5          if  $W[j, i]$  is not a substring of  $X$  then
6               $z \leftarrow z + 1$ 
7               $j \leftarrow i + 1$ 
8       $D(i) \leftarrow z$ 

```

Cuadro 2: Algoritmo de la función CALCULATED.[2]

2.2.2. Función de búsqueda inexacta recursiva InexRecur

Esta función es la encargada de realizar la búsqueda inexacta de forma recursiva. Tendrá en cuenta los diferentes tipos posibles de operaciones (match, sustituciones, deletiones e inserciones). Los parámetros de entrada de InexRecur serán el substring W , el número de errores z , y las variables de recursión i , k y l .

```

1  INEXRECUR( $W, i, z, k, l$ )
2      if  $z < D(i)$  then
3          return  $\emptyset$ 
4      if  $i < 0$  then
5          return  $\{[k, l]\}$ 
6           $I \leftarrow \emptyset$ 
7       $I \leftarrow I \cup \text{INEXRECUR}(W, i - 1, z - 1, k, l)$ 
8      for each  $b \in \{A, C, G, T\}$  do
9           $k \leftarrow C(b) + O(b, k - 1) + 1$ 
10          $l \leftarrow C(b) + O(b, l)$ 
11         if  $k \leq l$  then
12              $I \leftarrow I \cup \text{INEXRECUR}(W, i, z - 1, k, l)$ 
13             if  $b = W[i]$  then
14                  $I \leftarrow I \cup \text{INEXRECUR}(W, i - 1, z, k, l)$ 
15             else
16                  $I \leftarrow I \cup \text{INEXRECUR}(W, i - 1, z - 1, k, l)$ 
17     return  $I$ 

```

Cuadro 3: Algoritmo de la función INEXRECUR.[2]

Vemos pues que para cada posible tipo de operaciones el algoritmo llevará a cabo una acción concreta.

En el caso de que se haya producido una delección (borrado de un carácter o nucleótido en dicha posición) se llamará a la función recursiva disminuyendo una unidad la variable de posición y el número de errores (i, z). Si en cambio se ha observado una inserción (adición de un carácter o nucleótido en la secuencia) se llamará a la función `InexRecur` y se disminuirá en un tanto el número de errores z .

En el caso de que se haya producido un match (coincidencia) o una mutación de tipo sustitución se volverá a llamar a la función `InexRecur` pero disminuyendo en una unidad la variable i .

Y por último, si se ha producido una sustitución, se llamará a la función `InexRecur` disminuyendo en un tanto la variable de posición y el número de errores (i, z), como ocurría en el caso de una delección.

2.2.3. Función `InexSearch`

Esta función será la encargada de llamar a la función `InexRecur`, pasándole sus correspondientes parámetros de entrada W, i, z, k y l . Las variables de entrada a la función `InexSearch` serán por tanto el substring W y el umbral z . Podemos ver la lógica de esta función en la Figura 5.

<pre> 1 INEXACTSEARCH(W, z) 2 CALCULATED(W) 3 return INEXRECUR($W, W - 1, z, 1, X - 1$) </pre>

Cuadro 4: Algoritmo de la función `INEXSEARCH`. [2]

3. Materiales

Los lenguajes que hemos comparado son R, Python y C. La implementación en R es la misma a la usada en las prácticas de la asignatura. Con la implementación proporcionada hemos dividido dos versiones, una que dibuja la traza y otra que calcula el resultado de forma limpia. Todas las versiones del código escrito está disponible en el anejo.

El tiempo y la memoria la hemos medido de dos maneras. Primero hemos evaluado a la función únicamente usando herramientas del lenguaje. En Python hemos usado los paquetes `resource`[5] y `time`[6] para medir el tiempo y la memoria. En R hemos usado la utilidad `system.time`[7] para el tiempo y la librería `profmem`[8] para perfilar el uso de memoria. En C hemos medido el tiempo usando la librería `time.h`[9] y la memoria usando la herramienta `Valgrind`[10]. Después hemos medido el uso de recursos desde la línea de comandos, llamando a los programas como scripts usando funcionalidades propias de cada lenguaje. Con ello podemos evaluar el sobrecoste de cada una de las tres plataformas usadas. Para ello hemos empleado la herramienta `bench`[11] y el comando `time`[9].

4. Resultados

Hemos implementado las trazas de manera que son perfectamente idénticas en Python y en C . A continuación mostramos algunos ejemplos comparativos para verificar que todos los programas se comportan de forma igual. Posteriormente comparamos el uso de tiempo y memoria.

1		Mutation	i	z	k	l
2	Deletion	[l]	1	-1	0	6
3	Insertion	[g]	2	-1	1	2
4	Substitution	[g] -> [l]	1	-1	1	2
5	Insertion	[l]	2	-1	3	3
6	Match	[l]	1	0	3	3
7	Deletion	[o]	0	-1	3	3
8	Insertion	[o]	1	-1	5	5
9	Match	[o]	0	0	5	5
10	Deletion	[g]	-1	-1	5	5
11	Insertion	[g]	0	-1	1	1
12	Match	[g]	-1	0	1	1
13	----- {1,1} -----					
14	Insertion	[o]	2	-1	4	6
15	Substitution	[o] -> [l]	1	-1	4	6

Cuadro 5: Traza de INEXRECUR con X = “googol\$”, W = “gol”, z = 0 en C y Python

1		INEXRECUR	-	by XAVI GABRI AITANA ALFREDO	
2	-	D	[l]	1 -1 0 6	
3	-	I	[g]	2 -1 1 2	
4	-	S	[g -> l]	1 -1 1 2	
5	-	I	[l]	2 -1 3 3	
6	-	M	[l]	1 0 3 3	
7	-	-	D [o]	0 -1 3 3	
8	-	-	I [o]	1 -1 5 5	
9	-	-	M [o]	0 0 5 5	
10	-	-	- D [g]	-1 -1 5 5	
11	-	-	- I [g]	0 -1 1 1	
12	-	-	- M [g]	-1 0 1 1	
13	-	-	-	-	-> [c(1, 1)]
14	-	I	[o]	2 -1 4 6	
15	-	S	[o -> l]	1 -1 4 6	

Cuadro 6: Traza de INEXRECUR con X = “googol\$”, W = “gol”, z = 0 en R

1		Mutation	i	z	k	l
2	Deletion	[g]	2	-1	0	6
3	Insertion	[g]	3	-1	1	2
4	Match	[g]	2	0	1	2
5	Deletion	[o]	1	-1	1	2
6	Insertion	[o]	2	-1	4	4
7	Match	[o]	1	0	4	4
8	Deletion	[o]	0	-1	4	4
9	Insertion	[o]	1	-1	6	6
10	Match	[o]	0	0	6	6
11	Insertion	[l]	3	-1	3	3
12	Substitution	[l] -> [g]	2	-1	3	3
13	Insertion	[o]	3	-1	4	6
14	Substitution	[o] -> [g]	2	-1	4	6

Cuadro 7: Traza de INEXRECUR con X = “googol\$”, W = “goog”, z = 0 en C y Python

1		INEXRECUR	-	by XAVI GABRI AITANA ALFREDO	
2	-	D	[g]	2 -1 0 6	
3	-	I	[g]	3 -1 1 2	
4	-	M	[g]	2 0 1 2	
5	-	-	D [o]	1 -1 1 2	
6	-	-	I [o]	2 -1 4 4	
7	-	-	M [o]	1 0 4 4	
8	-	-	- D [o]	0 -1 4 4	
9	-	-	- I [o]	1 -1 6 6	
10	-	-	- M [o]	0 0 6 6	
11	-	I	[l]	3 -1 3 3	
12	-	S	[l -> g]	2 -1 3 3	
13	-	I	[o]	3 -1 4 6	
14	-	S	[o -> g]	2 -1 4 6	

Cuadro 8: Traza de INEXRECUR con X = “googol\$”, W = “goog”, z = 0 en R

1		Mutation	i	z	k	l
2	Deletion	[l]	2	0	1	6
3	Deletion	[o]	1	-1	1	6
4	Insertion	[g]	2	-1	1	2
5	Substitution	[g] -> [o]	1	-1	1	2
6	Insertion	[o]	2	-1	4	6
7	Match	[o]	1	0	4	6
8	Deletion	[o]	0	-1	4	6
9	Insertion	[g]	1	-1	1	2
10	Substitution	[g] -> [o]	0	-1	1	2
11	Insertion	[o]	1	-1	6	6
12	Match	[o]	0	0	6	6
13	Deletion	[g]	-1	-1	6	6
14	Insertion	[g]	0	-1	2	2
15	Match	[g]	-1	0	2	2
16	<hr/> {2,2} <hr/>					
17	Insertion	[g]	3	0	1	2
18	Substitution	[g] -> [l]	2	0	1	2
19	Deletion	[o]	1	-1	1	2
20	Insertion	[o]	2	-1	4	4
21	Match	[o]	1	0	4	4
22	Deletion	[o]	0	-1	4	4
23	Insertion	[o]	1	-1	6	6
24	Match	[o]	0	0	6	6
25	Deletion	[g]	-1	-1	6	6
26	Insertion	[g]	0	-1	2	2
27	Match	[g]	-1	0	2	2
28	<hr/> {2,2} <hr/>					
29	Insertion	[o]	3	0	4	6
30	Substitution	[o] -> [l]	2	0	4	6
31	Deletion	[o]	1	-1	4	6
32	Insertion	[g]	2	-1	1	2
33	Substitution	[g] -> [o]	1	-1	1	2
34	Insertion	[o]	2	-1	6	6
35	Match	[o]	1	0	6	6
36	Deletion	[o]	0	-1	6	6
37	Insertion	[g]	1	-1	2	2
38	Substitution	[g] -> [o]	0	-1	2	2

Cuadro 9: Traza de INEXRECUR con $X = \text{"googol\$"} , W = \text{"gool"} , z = 1$ en C y Python

[illegible]

Cuadro 10: Traza de INEXRECUR con $X = \text{"googol\$"}, W = \text{"gool"}, z = 1$ en R

Se puede comprobar que las tres trazas contienen la misma información. Procedemos ahora a la comparación de los tiempos medidos, primero las mediciones tomadas desde el propio programa.

```

1  inexrecur_time.c
2  12.34μs from 10000 iterations.
3
4  inexrecur_time.py
5  real      sys      user
6  268.89μs  0.37μs  268.16μs
7
8  inexrecur_time.R
9  user      system    elapsed
10 5422μs     18μs      5443μs

```

Cuadro 11: Tiempos de “CPU”.

En cada caso tomamos múltiples lecturas y hacemos la media. Como es de esperar la solución compilada de C es ~ 20 veces más rápida al script de Python y R parece ser ~ 450 veces más lento que C y ~ 20 veces más lento que Python . Para comparar el efecto del interpretador medimos ahora el tiempo ejecutando desde la línea de comandos. Llamamos a los scripts usando `#!/bin/env Rscript` y `#!/bin/env Python`.

```

1  bench ./inexrecur_clean ./inexrecur_clean.py ./inexrecur_clean.R
2  benchmarking bench/./inexrecur_clean
3  time                4.524 ms      (4.496 ms .. 4.551 ms)
4                      0.999 R2      (0.999 R2 .. 1.000 R2)
5  mean                4.522 ms      (4.499 ms .. 4.559 ms)
6  std dev             89.83 μs      (58.07 μs .. 133.1 μs)
7
8  benchmarking bench/./inexrecur_clean.py
9  time                39.26 ms      (39.12 ms .. 39.40 ms)
10                      1.000 R2      (1.000 R2 .. 1.000 R2)
11  mean                39.30 ms      (39.18 ms .. 39.45 ms)
12  std dev             266.5 μs      (177.4 μs .. 388.7 μs)
13
14  benchmarking bench/./inexrecur_clean.R
15  time                278.5 ms      (273.6 ms .. 285.0 ms)
16                      1.000 R2      (1.000 R2 .. 1.000 R2)
17  mean                280.6 ms      (279.1 ms .. 281.9 ms)
18  std dev             1.714 ms      (1.027 ms .. 2.517 ms)
19  variance introduced by outliers: 16% (moderately inflated)

```

Cuadro 12: Tiempo de “pared” según la utilidad *bench*.

En este caso C es ~ 10 veces más rápido a Python y ~ 70 veces más rápido que R , que continua siendo ~ 7 veces más lento que Python .

```

1  inexrecur_clean.c
2  ==81469==      in use at exit: 18,588 bytes in 164 blocks
3  ==81469==      total heap usage: 191 allocs , 27 frees , 24,894 bytes allocated
4
5  inexrecur_mem.py
6  6,631,424 bytes
7
8  inexrecur_mem.R
9  4,932,584 bytes

```

Cuadro 13: Memoria usada medido desde el script.

La memoria usada es comparable en R y Python cuando la medimos sin tener en cuenta el entorno de ejecución, C utiliza entorno a 200 veces menos memoria.

```

1  time ./inexrecur_clean
2  (5,5)
3  (1,1)
4  ./inexrecur_clean
5  0.00s  user 0.00s system 44% cpu 0.004 total
6  max memory:          676 kB
7
8  time ./inexrecur_clean.py
9  (1, 1)
10 (5, 5)
11 ./inexrecur_clean.py
12 0.03s  user 0.01s system 82% cpu 0.039 total
13 max memory:          6272 kB
14
15 time ./inexrecur_clean.R
16 [[1]]
17 [1] 5 5
18
19 [[2]]
20 [1] 1 1
21
22 ./inexrecur_clean.R
23 0.23s  user 0.04s system 97% cpu 0.277 total
24 max memory:          67476 kB

```

Cuadro 14: Memoria usada medido desde fuera del script.

Cuando consideramos el entorno de ejecución interpretado al completo C se mantiene a la cabeza con la sorpresa de que R emplea ~10 veces más memoria que Python en este ejemplo en particular.

5. Discusión

Como se ha visto en el apartado de Resultados, hemos implementado el algoritmo BWA en los lenguajes R, Python y C; hemos utilizado distintos criterios para evaluar el rendimiento de éstos y realizar una comparación de su funcionamiento. Dichos criterios son el tiempo de procesamiento y la memoria usada. En el caso del tiempo, centrándonos en los dos lenguajes que queremos comparar en este artículo, el que ofrece un tiempo menor medido desde el propio programa a la hora de alinear las secuencias es Python, mientras que R resulta ser más lento en comparación con éste. Por otro lado, evaluando la memoria, Python ocupa aproximadamente 10 veces menos memoria que R midiendo desde fuera del script.

Python es un lenguaje con mensaje interactivo, permite la escritura dinámica y administra la memoria automáticamente. Tiene una sintaxis muy clara e intuitiva, y utiliza márgenes para hacer los bloques de declaraciones. Además, se trata de un software libre por lo que se puede acceder a él fácilmente y desde todos los sistemas operativos. [12]

R se trata de un lenguaje de programación con enfoque de análisis estadístico, no se suele usar para implementar programas largos como se hace con otros lenguajes. Permite combinar programación imperativa con funcional, además de incluir programación orientada a objetos. Asimismo, es uno de los lenguajes más utilizados en investigación científica. [13, 14]

Así pues, basándonos en los ejemplos que hemos ejecutado, parece ser que en Python es más sencillo implementar el algoritmo InexRecur (ver Anexo) y además presenta un mejor valor en los criterios evaluados anteriormente. Por lo tanto, dado que alinear secuencias es un procedimiento que ya de por sí consume mucho tiempo, podríamos decir que a la hora de mapear una secuencia corta con un genoma de referencia sería más rápido y eficaz utilizar el lenguaje Python [15]. Sin embargo, hay que tener en cuenta que como este algoritmo es computacionalmente costoso no será conveniente utilizarlo para alinear secuencias cortas con, por ejemplo, el genoma humano completo, ya que daría unos tiempos de procedimiento demasiado elevados [16]. De la misma forma, también se comenta este mismo inconveniente en el artículo de Heng Li et al., 2009 [2] donde se expone que el rendimiento del alineador BWA se reduce cuando son lecturas largas, por lo que aconsejan que a la hora de trabajar con secuencias largas se divida la lectura en varios fragmentos cortos, alinearlos por separado y posteriormente unir las alineaciones parciales en una global.

6. Conclusiones

Hemos propuesto una implementación del algoritmo InexRecur en Python para alineamientos inexactos de fragmentos de secuencias con secuencias largas. Lo implementamos de forma correcta y comprobamos su funcionamiento con varios ejemplos con la traza que se nos ha proporcionado. Se comprueba que la programación de este algoritmo en Python proporciona un gran potencial para disminuir el tiempo de procesamiento y la memoria total utilizada, además de ofrecer una sintaxis más clara en comparación con la implementación en R.

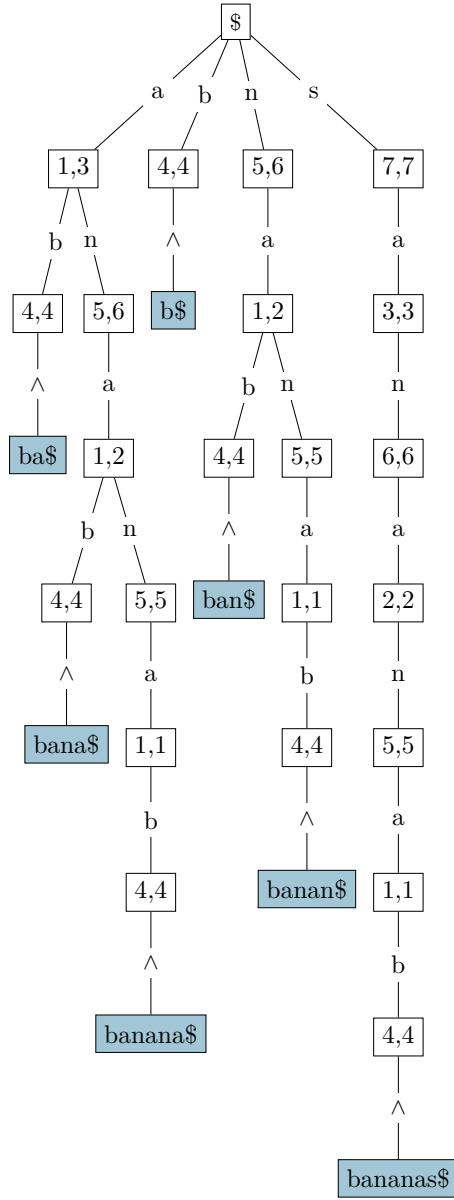


Figura 4: Árbol de prefijos con intervalos SA de la palabra “bananas\$”.

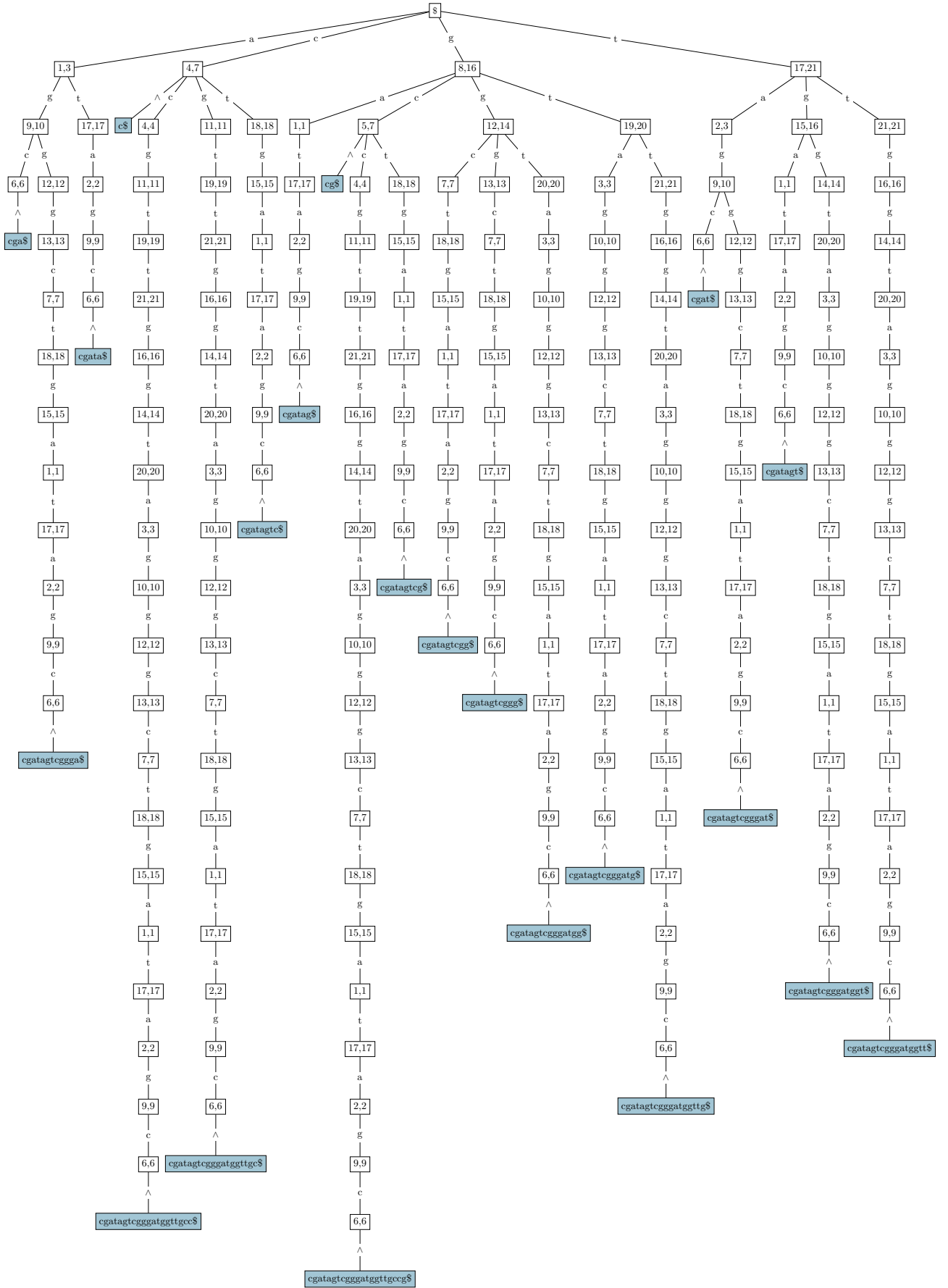


Figura 5: Árbol de prefijos con intervalos SA de la palabra “cgatagtcgggatggttgcg\$”.

7. Anejo

Listing 1: Python inexrecur_clean.py

```
1 #!/usr/bin/env python
2  # Started Wed 25 Mar 00:56:18 2020 by nahnero. #
3  def CalculatedD (W):
4      D = [0]*len (W)
5      z, j = 0, 0
6      for i in range (len (W)):
7          if W[j:i] not in X:
8              z += 1
9              j = i + 1
10         D[i] = z
11     return D
12
13 def InexRecur (W, i, z, k, l):
14     if i < 0:
15         if z < 0:
16             return set ([])
17         else:
18             return set ([ (k, l) ])
19     if z < D[i]:
20         return set ([])
21     I = set ([])
22     I = I.union (InexRecur (W, i-1, z-1, k, l))
23     for bi in dic:
24         ki = C (bi) + 0 (bi, k-1) + 1
25         li = C (bi) + 0 (bi, l)
26         if ki <= li:
27             I = set (I.union (InexRecur (W, i, z-1, ki, li)))
28             if bi == W[i]:
29                 I = set (I.union (InexRecur (W, i-1, z, ki, li)))
30             else:
31                 I = set (I.union (InexRecur (W, i-1, z-1, ki, li)))
32     return I
33
34 def InexSearch (W, z):
35     return InexRecur (
36         W = W,
37         i = len (W)-1,
38         z = z,
39         k = 0,
40         l = len (X)-1)
41 X = 'googol$';
42 W = 'lol'
43 dic = ''.join (sorted (set (X.replace ('$ ', ' '))))
44
45 BWT = [(X*2)[i:i+len (X)], i] for i in range (len (X))]
46 BWT.sort ()
47 S, B = [], ''
48 for i in BWT: S.append (i[1]); B += (i[0][-1])
49 C = lambda a: ''.join (sorted (X.replace ('$ ', ' ')))[0:-2].find (a))
50 O = lambda a, i: B[0:i+1].count (a)
51
52 D = CalculatedD (W)
53
54 IS = InexSearch (W, 1)
55 for i in IS: print (i)
```

Listing 2: R inexrecur_clean.R

```

1  #!/usr/bin/env Rscript
2  C <- function(a) {
3    components <- gsub("[[:punct:]]", "", X)
4    components <- strsplit(components, "")
5    components_sorted <- sort(components[[1]])
6    idx <- min(which(components_sorted == a))
7    C_num <- idx-1
8    return(C_num)
9  }
10 O <- function(a,i) {
11   i<-i+1
12   coincidences <- grep(a,B)
13   counter <- i>=coincidences
14   O_num <- length(counter[counter==TRUE])
15   return(O_num)
16 }
17
18 INEXRECUR <- function(W,i,z,k,l,sep) {
19   if (i<0){
20     if (z<0) {return({})}
21     else {
22       result <- list(c(k,l))
23       return(result)
24     }
25   if (z<D[i+1]) {return({})}
26   I <- {}
27   I <- c (I, INEXRECUR(W,i-1,z-1,k,l,sep));
28   for (bi in b){
29     ki <- C(bi) + O(bi,k-1) + 1 #Precalculate C and O function
30     li <- C(bi) + O(bi,l) #Precalculate C and O function
31     if (ki<=li){
32       I <- c (I, INEXRECUR(W,i,z-1,ki,li,sep));
33       if (bi==W[[1]][i+1]){
34         I <- c (I, INEXRECUR(W,i-1,z,ki,li,sep))
35       }else{
36         I <- c (I, INEXRECUR(W,i-1,z-1,ki,li,sep))
37       }
38     }
39   }
40   return(I)
41 } #end INEXRECUR
42
43 X <- "googol$"
44 SA_pre <- c()
45 SA_pre[1] <- X
46 for (i in 2:nchar(X)){
47   SA_pre[i] <- paste(substring(X,i,nchar(X)),
48                     substring(X,1,i-1),sep='',collapse=NULL)
49 }
50 SA_list<-sort(SA_pre,index.return=TRUE)
51 SA<-SA_list$x
52 S<-SA_list$ix
53 B<-paste(substring(SA,nchar(X),nchar(X)),sep='',collapse=NULL)
54 b <- c("g","l","o")
55 D <- c(0,0,1)
56 W = list(c("l","o","l"))
57 i = 2
58 z = 1
59 a = INEXRECUR(W,i,z,0,6,"")
60 print (a)

```

Listing 3: C inexrecur_clean.c

```

1  int main (void){
2      char X[] = "googol$";
3      char W[] = "lol";
4      char dict[] = "glo";
5      setup (X, W, dict);
6      return 0;
7  }
8
9
10 void setup (char* X1, char* W1, char* dic1){
11
12     int i, length; i = length = 0;
13     X = X1;
14     /* reverse */
15     W = W1;
16     while (W1[++i]);
17     Wlen = i;
18     dic = dic1;
19
20     while (X[++length]);
21
22     D = (int*) calloc (length, sizeof (int));
23     BWT = (char**)calloc (length, sizeof (char*));
24     Xs = (char*) calloc (length, sizeof (char));
25     B = (char*) calloc (length, sizeof (char));
26     X2 = (char*) calloc (2*length+1, sizeof (char));
27     for (i = 0; i != length; i++)
28         BWT[i] = (char*)malloc (length*sizeof (char)+1);
29
30     sprintf (X2, "%s%s", X, X);
31
32     for (i = 0; i != length; i++)
33         sprintf (BWT[i], "%.s", (int)length, X2+i);
34
35
36     sort (BWT, length);
37     CalculatedD (W1);
38
39     for (i = 0; i != length; i++) B [i] = BWT[i][length-1];
40     for (i = 0; i != length; i++) Xs[i] = BWT[i][0];
41
42     i = 0; while (X[++i]);
43     InexRecur (W1, Wlen - 1, 1, 0, i - 1);
44
45     free (Xs);
46     free (X2);
47     free (D);
48     for (i = 0; i != length ; i++) free (BWT[i]);
49     free (BWT);
50 }
51
52 static int myCompare(const void* a, const void* b){
53     return strcmp (*(const char**)a, *(const char**)b);
54 }
55
56 void sort(char* arr[], int n){
57     qsort (arr, n, sizeof(char*), myCompare);
58 }

```

Listing 4: C inexrecur_clean.c (cont.)

```

59
60 int C (char a){
61     int i = 0;
62     while (X[++i]) if (Xs[i] == a) return i - 1;
63     return 0;
64 }
65
66 int O (char a, int i){
67     int acc = 0, j;
68     for (j = 0; j <= i; j++) {
69         if (B[j] == a) acc++;
70     }
71     return acc;
72 }
73
74 /* Ignacio Amat */
75
76 void CalculatedD (char* W){
77     short z = 0, j = 0, i;
78     char buf[BUFSIZ] = {'\0'};
79     extern int Wlen;
80     for (i = 0; i != Wlen; i++){
81         sprintf (buf, "%.s", i, W+j);
82         if (!strstr (X, buf)){
83             z++;
84             j = i + 1;
85         }
86         D[i] = z;
87     }
88 }
89
90 void InexRecur (char* W, int i, int z, int k, int l){
91     if (i < 0){
92         if (z < 0){
93             return;
94         } else {
95             printf ("%d,%d\n", k, l);
96             return;
97         }
98     }
99     if (z < D[i]){return;}
100     InexRecur (W, i-1, z-1, k, l);
101     int ki, li, m;
102     for (m = 0; m != 3; m++){
103         ki = C (dic [m]) + O (dic [m], k-1) + 1;
104         li = C (dic [m]) + O (dic [m], l);
105         if (ki <= li){
106             InexRecur (W, i, z-1, ki, li);
107             if (dic [m] == W[i]){
108                 InexRecur (W, i-1, z, ki, li);
109             } else {
110                 InexRecur (W, i-1, z-1, ki, li);
111             }
112         }
113     }
114     return;
115 }

```

Referencias

- [1] D. Giménez, “Implementación y análisis de algoritmos de alineación para datos de next generation sequencing (ngs),” Ph.D. dissertation, Universidad Autónoma de Madrid (UAM), 2016.
- [2] H. Li and R. Durbin, *Fast and accurate short read alignment with Burrows-Wheeler transform.*, 2009. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btp324>
- [3] J. M. García Gómez, *Bioinformatics Burrows-Wheeler Alignment (BWA)*. Apuntes de la asignatura de Bioinformática, Grado en Ingeniería Biomédica, UPV, 2020. [Online]. Available: https://poliformat.upv.es/access/content/group/GRA_13024_2019/JuanMGarcia-Gomez/Semana%204/BWA.pdf
- [4] J. G. García Pardo, A. Torregrosa Lloret, J. De la Torre Costa, and A. Pascual Belda, *InexRecur*. Apuntes de la asignatura de Bioinformática, Grado en Ingeniería Biomédica, UPV, 2020. [Online]. Available: https://poliformat.upv.es/access/content/group/GRA_13024_2019/JuanMGarcia-Gomez/Semana%204/INEXRECUR.pdf
- [5] “resource - resource usage information.” [Online]. Available: <https://docs.python.org/3/library/resource.html>
- [6] “time - time access and conversions.” [Online]. Available: <https://docs.python.org/3/library/time.html?highlight=time#module-time>
- [7] “system.time.” [Online]. Available: <https://www.rdocumentation.org/packages/base/versions/3.6.2/topics/system.time>
- [8] H. Bengtsson, “Simple memory profiling for r [r package profmem version 0.5.0].” [Online]. Available: <https://cran.r-project.org/web/packages/profmem/index.html>
- [9] “time.h,” <https://pubs.opengroup.org/onlinepubs/007908799/xsh/time.h.html>, (Accessed on 05/02/2020).
- [10] “Valgrind home,” <https://valgrind.org/>, (Accessed on 05/02/2020).
- [11] “Github - gabriel439/bench: Command-line benchmark tool,” <https://github.com/Gabriel439/bench>, (Accessed on 05/02/2020).
- [12] H. Fangohr, *A Comparison of C, MATLAB, and Python as Teaching Languages in Engineering*, 2004. [Online]. Available: https://link.springer.com/content/pdf/10.1007%2F978-3-540-25944-2_157.pdf
- [13] C. J. Gil Bellosta, *R para profesionales de los datos: una introducción*, 2018. [Online]. Available: https://www.datanalytics.com/libro_r/programacion-en-r.html
- [14] C. Kleiber and A. Zeileis, *Applied Econometrics with R*, 1st ed. Springer Science+Business Media, 2008.
- [15] J. M. Abuín, J. C. Pichel, T. F. Pena, and J. Amigo, *BigBWA: approaching the Burrows-Wheeler aligner to Big Data technologies*. International Society for Computational Biology, 2015. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btv506>

- [16] E. J. Houtgast, V.-M. Sima, K. Bertels, and Z. Al-Ars, “An FPGA-based systolic array to accelerate the BWA-MEM genomic mapping algorithm,” in *2015 International Conference on Embedded Computer Systems: Architectures, Modeling, and Simulation (SAMOS)*. IEEE, Jul. 2015. [Online]. Available: <https://doi.org/10.1109/samos.2015.7363679>