

String Matching usando Transformada Discreta de Fourier

Notación

Considerando el problema de matching exacto de cadenas de caracteres, sean:

- t : Un texto
- p : Un patrón a buscar dentro del texto
- $t[n]$: Una representación numérica del n -ésimo caracter de t
- $p[n]$: Una representación numérica del n -ésimo caracter de p
- N : La longitud del texto t
- M : La longitud del patrón p
- $((x))_a$: La operación $x \bmod a$

Planteo Formal

El problema se puede escribir como:

$$\sum_{k=0}^{M-1} |t[k+i] - p[k]| = 0 \Leftrightarrow \text{Hay un match que comienza en la posición } i.$$

O bien, por conveniencia, calculamos la convergencia cuadrática (notar que no estamos elevando ambos miembros al cuadrado, sino cada término):

$$\sum_{k=0}^{M-1} (t[k+i] - p[k])^2 = 0 \Leftrightarrow \text{Hay un match que comienza en la posición } i.$$

Desarrollando cada binomio de la sumatoria:

$$\sum_{k=0}^{M-1} (t[k+i]^2 - 2t[k+i]p[k] + p[k]^2) = 0$$

Multiplicando punto a punto por el arreglo $t[i+k]$:

$$\sum_{k=0}^{M-1} (t[k+i]^3 - 2t[k+i]^2p[k] + t[k+i]p[k]^2) = 0$$

Y haciendo lo propio por el arreglo $p[k]$:

$$\sum_{k=0}^{M-1} (t[k+i]^3p[k] - 2t[k+i]^2p[k]^2 + t[k+i]p[k]^3) = 0$$

Extendiendo con ceros el arreglo p , puedo cambiar los límites de la sumatoria:

$$\sum_{k=0}^{N-1} \left(t[((k+i))_N]^3 p[((k))_N] - 2t[((k+i))_N]^2 p[((k))_N]^2 + t[((k+i))_N] p[((k))_N]^3 \right) = 0$$

Por asociatividad de la suma:

$$\sum_{k=0}^{N-1} t[((k+i))_N]^3 p[((k))_N] - 2 \sum_{k=0}^{N-1} t[((k+i))_N]^2 p[((k))_N]^2 + \sum_{k=0}^{N-1} t[((k+i))_N] p[((k))_N]^3 = 0$$

Aplicando la definición de convolución circular:

$$t[i]^3 \odot p[M-1-i] - 2t[i]^2 \odot p[M-1-i]^2 + t[i] \odot p[M-1-i]^3 = 0$$

Sean:

- $T_i[k]$: La DFT de $t[n]^i$.
- $P_i[k]$: La DFT de $p[n]^i$.

Y las siguientes propiedades de la transformada discreta de Fourier:

- $DFT_N \{ x[((n-a))_N] \} = e^{-j(2\pi/N) a k} X[k]$
- $DFT_N \{ x^*[((-n))_N] \} = X^*[k]$

Puedo transformar ambos miembros de la ecuación como:

$$T_3[k] P_1^*[k] e^{j(2\pi/N) (M-1) k} - 2 T_2[k] P_2^*[k] e^{j(2\pi/N) (M-1) k} + T_1[k] P_3^*[k] e^{j(2\pi/N) (M-1) k} = 0$$

Aplicando los productos puntuales en cada término:

$$R_1[k] e^{j(2\pi/N) (M-1) k} - 2 R_2[k] e^{j(2\pi/N) (M-1) k} + R_3[k] e^{j(2\pi/N) (M-1) k} = 0$$

Y aplicando la transformada inversa:

$$r_1[((i+M-1))_N] - 2 r_2[((i+M-1))_N] + r_3[((i+M-1))_N] = 0$$

Finalmente, si se cumple esta ecuación, hay un match que comienza en la posición i de t .

Además, se observan las siguientes igualdades:

$$\begin{aligned} r_1[((i+M-1))_N] &= t[i]^3 \odot p[M-1-i] \\ r_2[((i+M-1))_N] &= t[i]^2 \odot p[M-1-i]^2 \\ r_3[((i+M-1))_N] &= t[i] \odot p[M-1-i]^3 \end{aligned}$$

Algoritmo

A partir del planteo anterior y utilizando como primitiva una implementación de la transformada rápida de Fourier (con orden algorítmico $O(N \cdot \log(N))$), podemos escribir el siguiente algoritmo:

```
# Sea la Transformada Rápida de Fourier (Fast Fourier Transform):
Arreglo<N> FFT(Arreglo<N>)

# Se instancian los arreglos
t1, t2, t3, p1, p2, p3 = Arreglos de largo N inicializado con ceros

# Representación numérica del texto
for i in 0..N-1:
    t1[i] := Representación numérica del i-ésimo caracter del texto
    t2[i] := t1[i] * t1[i]
    t3[i] := t2[i] * t1[i]

# Representación numérica del patrón invirtiendo el dominio
for i in 0..M-1:
    p1[M-i-1] := Representación i-ésimo caracter del pattern
    p2[M-i-1] := p1[M-i-1] * p1[M-i-1]
    p3[M-i-1] := p2[M-i-1] * p1[M-i-1]

# Cálculo de las transformadas discretas de Fourier del texto
T1 := FFT(t1)
T2 := FFT(t2)
T3 := FFT(t3)

# Cálculo de las transformadas discretas de Fourier del patrón
P1 := FFT(p1)
P2 := FFT(p2)
P3 := FFT(p3)

# Productos punto a punto de las transformadas
R1 := T3 * P1
R2 := T2 * P2
R3 := T1 * P3

# Cálculo de las transformadas inversas de los productos
r1 := IFFT(R1)
r2 := IFFT(R2)
r3 := IFFT(R3)

for i in 0..N-M:
    if r1[M-1-i] - 2 * r2[M-1-i] + r3[M-1-i] == 0:
        Mostrar match en la posición i
```

Orden del Algoritmo

Como la **FFT** tiene orden $O(N \log(N))$, y el resto de las operaciones del algoritmo son de orden lineal, se conservará el orden de la **FFT**.

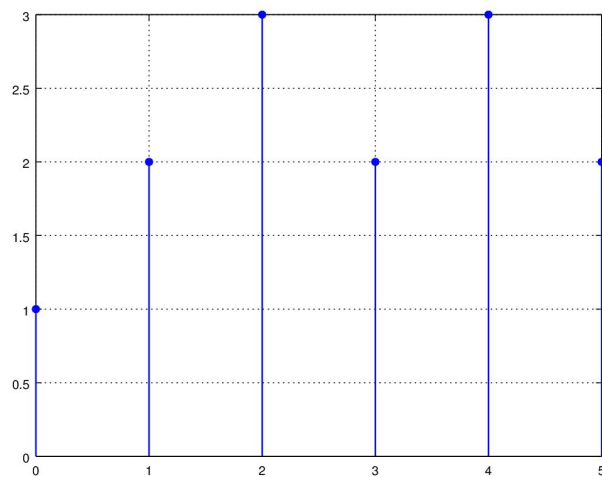
Ejemplo gráfico

Buscaremos la subcadena 'ANA' dentro de la cadena 'BANANA'.

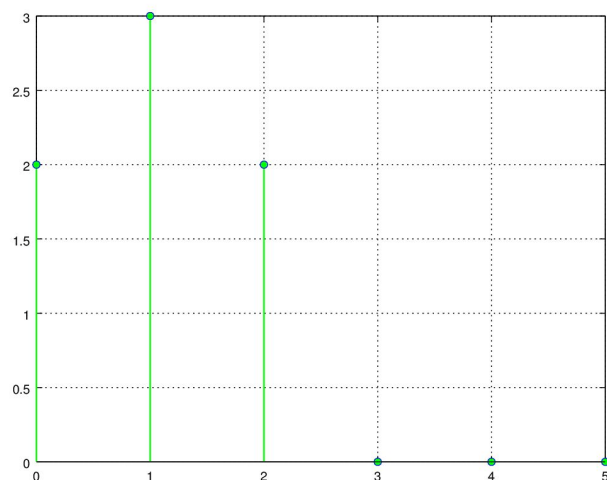
Tenemos entonces:

- $t = \{1, 2, 3, 2, 3, 2\}$
- $p = \{2, 3, 2, 0, 0, 0\}$
- El patrón invertido será también $p' = \{2, 3, 2, 0, 0, 0\}$
- $N = 6$
- $M = 3$

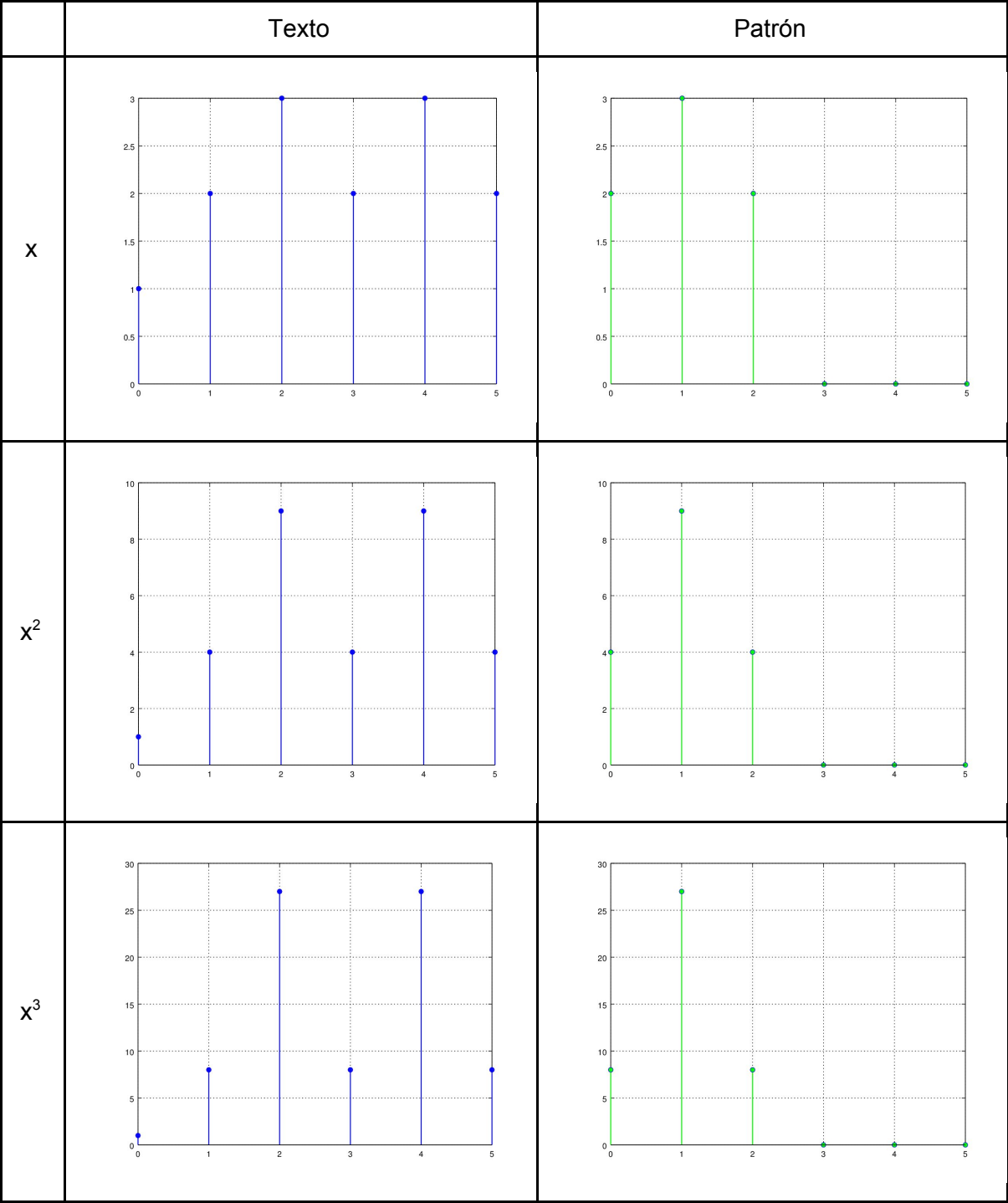
El texto representado como señal se verá como:



Y el patrón (ya invertido) como:



Las señales elevadas punto a punto al cuadrado y al cubo:

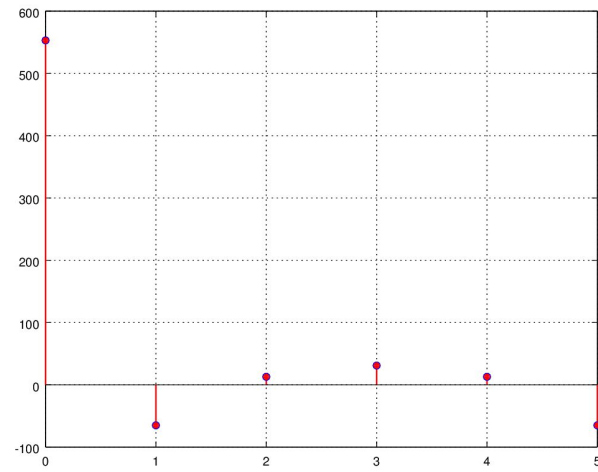


Sus transformadas de Fourier discretas (tomo parte real para graficar en dos dimensiones):

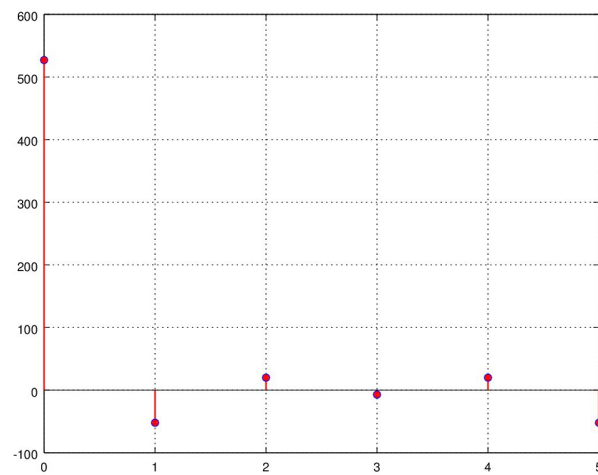
	Texto	Patrón
X_1	<p>Discrete Fourier Transform plot for X_1. The x-axis ranges from 0 to 5, and the y-axis ranges from -5 to 15. The plot shows a main peak at $k=0$ with a value of approximately 13.5, and smaller side lobes at $k=1, 2, 3, 4, 5$ with values of approximately -2.5, -2.5, 1.0, -2.5, and -2.5 respectively.</p>	<p>Discrete Fourier Transform plot for X_1. The x-axis ranges from 0 to 5, and the y-axis ranges from -2 to 8. The plot shows a main peak at $k=0$ with a value of approximately 7.5, and smaller side lobes at $k=1, 2, 3, 4, 5$ with values of approximately 2.5, -0.5, 1.0, -0.5, and 2.5 respectively.</p>
X_2	<p>Discrete Fourier Transform plot for X_2. The x-axis ranges from 0 to 5, and the y-axis ranges from -10 to 40. The plot shows a main peak at $k=0$ with a value of approximately 32, and smaller side lobes at $k=1, 2, 3, 4, 5$ with values of approximately -8, -8, 7, -8, and -8 respectively.</p>	<p>Discrete Fourier Transform plot for X_2. The x-axis ranges from 0 to 5, and the y-axis ranges from -5 to 20. The plot shows a main peak at $k=0$ with a value of approximately 18, and smaller side lobes at $k=1, 2, 3, 4, 5$ with values of approximately 6.5, -2.5, -1.0, -2.5, and 6.5 respectively.</p>
X_3	<p>Discrete Fourier Transform plot for X_3. The x-axis ranges from 0 to 5, and the y-axis ranges from -40 to 80. The plot shows a main peak at $k=0$ with a value of approximately 80, and smaller side lobes at $k=1, 2, 3, 4, 5$ with values of approximately -25, -25, 30, -25, and -25 respectively.</p>	<p>Discrete Fourier Transform plot for X_3. The x-axis ranges from 0 to 5, and the y-axis ranges from -20 to 50. The plot shows a main peak at $k=0$ with a value of approximately 45, and smaller side lobes at $k=1, 2, 3, 4, 5$ with values of approximately 18, -10, -12, -10, and 18 respectively.</p>

Haciendo el producto punto a punto de las señales:

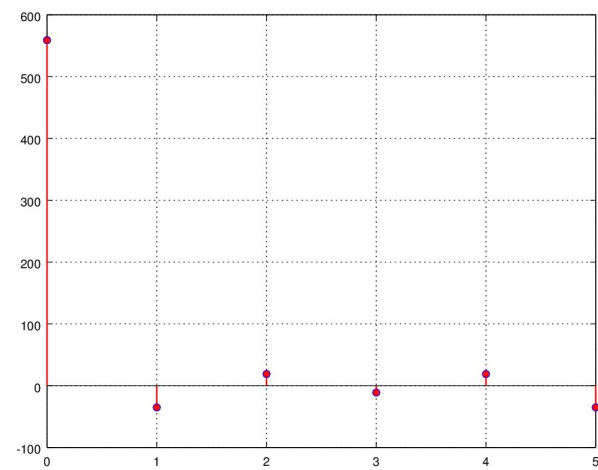
$$R_1 = T_3 P_1$$



$$R_2 = T_2 P_2$$

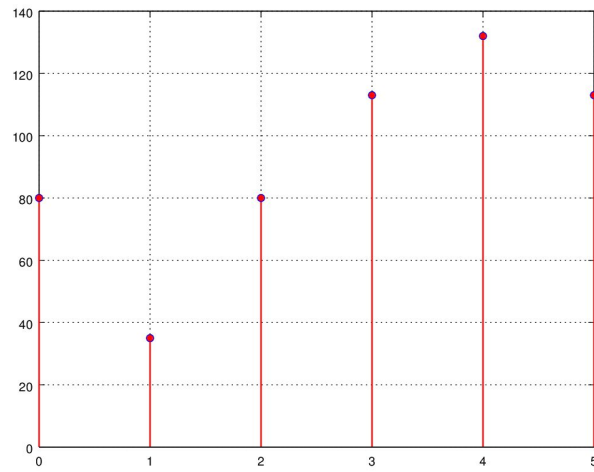


$$R_3 = T_1 P_3$$

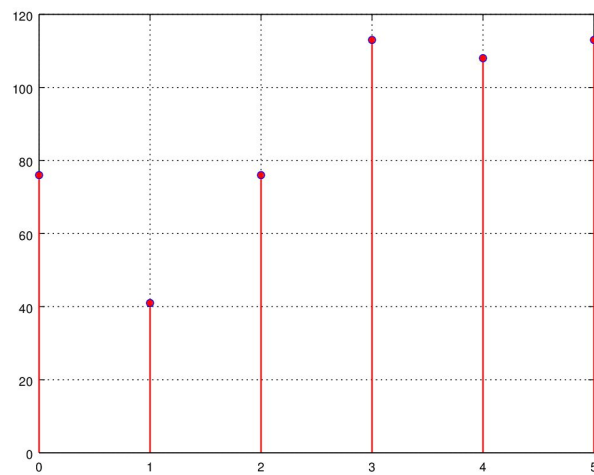


Aplicando la transformada discreta inversa de Fourier:

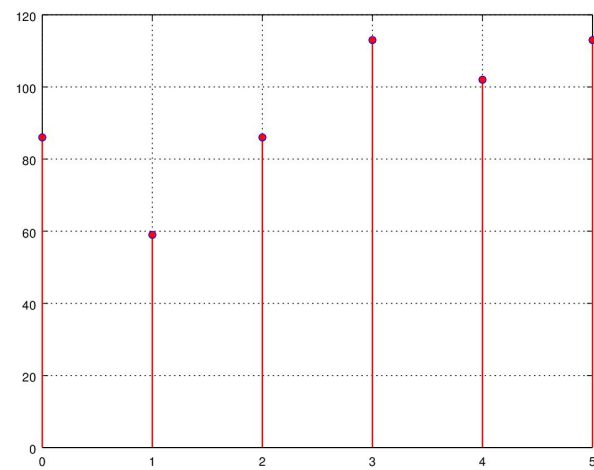
$$r_1 = \text{IDFT}(R_1)$$



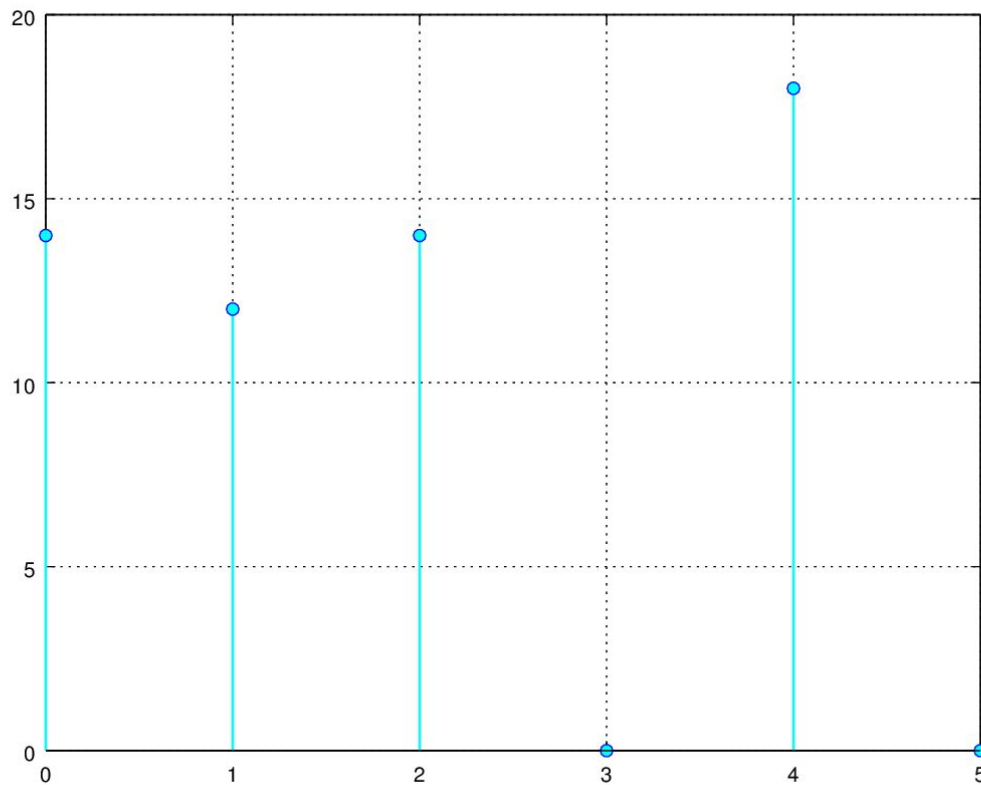
$$r_2 = \text{IDFT}(R_2)$$



$$r_3 = \text{IDFT}(R_3)$$



Y por último, aplicando la combinación lineal $r_1[i] - 2 r_2[i] + r_3[i]$ punto a punto:



Vemos que hay ceros en $i=3$ e $i=5$, lo cual quiere decir que los matches terminan en las posiciones 3 y 5 del texto (empiezan en $n=3+M-1=1$ y $n=5+M-1=3$).

Adaptaciones del Algoritmo

- Si buscamos sub-cadenas, sólo tenemos en cuenta las posiciones entre $M-1$ y $N-1$ de la señal resultante, las anteriores son útiles para buscar sub-cadenas en el texto visto de manera circular (encontrar un match desde la posición -2 hasta la 3, por ejemplo).
- Si queremos buscar con “*don't cares*” (cualquier caracter), simplemente tenemos que poner un cero en la representación del caracter correspondiente en la representación numérica del patrón.