

String Matching

Cristian López Del Alamo ¹,
Milton Condori Fernández ²
Jose Luis Sotomayor Malqui ²

Universidad Nacional de San Agustín
Escuela Profesional de Ciencia de la Computación

¹Universidad La Salle

²Universidad Nacional de San Agustín

Índice

- 1 Introducción
 - Motivación
 - String Matching
 - Formalidades
- 2 Algoritmos
- 3 Analisis
 - Algoritmo obvio
 - Algoritmo de Rabin Karp
 - Algoritmo del Automata Finito
 - Algoritmo de Knuth Morris Pratt

Motivación

- A menudo los programas de edición de texto tienen que encontrar todas las ocurrencias de un patrón brindado por el usuario en dicho texto.

Motivación

- A menudo los programas de edición de texto tienen que encontrar todas las ocurrencias de un patrón brindado por el usuario en dicho texto.
- Procesamiento de cadenas de ADN requiere el reconocimiento de patrones de secuencias específicas de ADNs.

Motivación

- A menudo los programas de edición de texto tienen que encontrar todas las ocurrencias de un patrón brindado por el usuario en dicho texto.
- Procesamiento de cadenas de ADN requiere el reconocimiento de patrones de secuencias específicas de ADNs.
- Buscadores de internet necesitan encontrar palabras específicas para poder brindar las paginas web correspondientes ante una consulta.

String Matching

Definición

Algoritmos que solucionan el problema de reconocer y ubicar patrones de cadenas dentro de un texto son denominados algoritmos de coincidencia(String Matching Algorithms)

Notación

- Texto $\rightarrow T[1..n]$

Notación

- Texto $\rightarrow T[1..n]$
- Patrón $\rightarrow P[1..m]$

Notación

- Texto $\rightarrow T[1..n]$
- Patrón $\rightarrow P[1..m]$
- Se dice que existe un matching entre P y T si: Existe un desplazamiento:

$$s, 0 \leq s \leq n-m,$$

Notación

- Texto $\rightarrow T[1..n]$
- Patrón $\rightarrow P[1..m]$
- Se dice que existe un matching entre P y T si: Existe un desplazamiento:

$$s, 0 \leq s \leq n-m,$$

tal que:

$$T[s+1..m] = P[1..m]$$

Notación

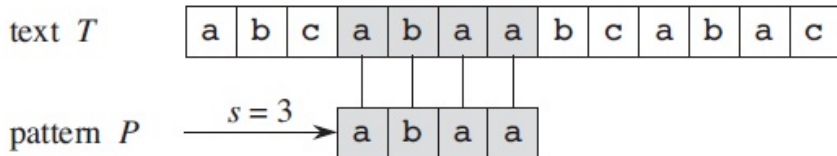
- Texto $\rightarrow T[1..n]$
- Patrón $\rightarrow P[1..m]$
- Se dice que existe un matching entre P y T si: Existe un desplazamiento:

$$s, 0 \leq s \leq n-m,$$

tal que:

$$T[s+1..m] = P[1..m]$$

String Matching



Notación

- “w es prefijo de x” ($w \sqsubset x$) si $x = w.y$ ($y \in \Sigma^*$).

Si $w \sqsubset x$ luego, $|w| \leq |x|$

Ejemplo:

$ab \sqsubset abcac$

Notación

- “w es prefijo de x” ($w \sqsubset x$) si $x = w.y$ ($y \in \Sigma^*$).

Si $w \sqsubset x$ luego, $|w| \leq |x|$

Ejemplo:

$ab \sqsubset abcac$

- “w es sufijo de x” ($w \sqsupset x$) si $x = y.w$

Ejemplo:

$cab \sqsupset abcab$

Notación

- “w es prefijo de x” ($w \sqsubseteq x$) si $x = w.y$ ($y \in \Sigma^*$).

Si $w \sqsubseteq x$ luego, $|w| \leq |x|$

Ejemplo:

$$ab \sqsubseteq abcac$$

- “w es sufijo de x” ($w \sqsubseteq x$) si $x = y.w$

Ejemplo:

$\text{cab} \sqsubset \text{abcab}$

- La notación P_k es equivalente a $P[1..k]$, representa a los k caracteres prefijos de P

Notación

- Ejemplo:

$$ab \sqsubset abcac$$

- Ejemplo:

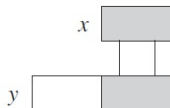
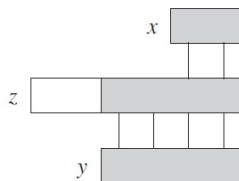
$\text{cab} \sqsubset \text{abcab}$

- La notación P_k es equivalente a $P[1..k]$, representa a los k caracteres prefijos de P
- T_k representa a los k caracteres prefijo de T

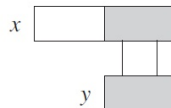
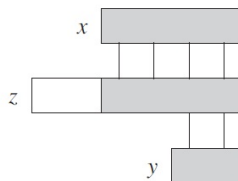
Notación

- x, y, z : cadenas tal que :

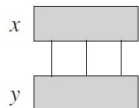
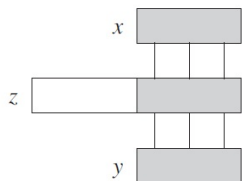
$x \sqsupset z$ y $y \sqsupset z$. Si $|x| < |y|$, luego $x \sqsubset y$. Si $|x| \geq |y|$, luego $y \sqsubset x$. Si $|x| = |y|$, luego $x = y$



(a)



(b)



(c)

Algunas Soluciones

Algorithm	Preprocessing time	Matching time
Naive	0	$O((n - m + 1)m)$
Rabin-Karp	$\Theta(m)$	$O((n - m + 1)m)$
Finite automaton	$O(m \Sigma)$	$\Theta(n)$
Knuth-Morris-Pratt	$\Theta(m)$	$\Theta(n)$

Algoritmo obvio

Este algoritmo es el que uno piensa naturalmente, posee un costo de computacional de $O((n - m + 1).m)$ en el peor de los casos.

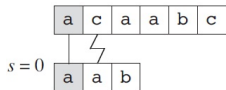
Algoritmo obvio

Este algoritmo es el que uno piensa naturalmente, posee un costo de computacional de $O((n - m + 1).m)$ en el peor de los casos.

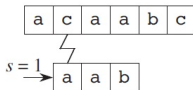
NAIVE-STRING-MATCHER(T, P)

```
1   $n = T.length$ 
2   $m = P.length$ 
3  for  $s = 0$  to  $n - m$ 
4      if  $P[1..m] == T[s + 1..s + m]$ 
5          print "Pattern occurs with shift"  $s$ 
```

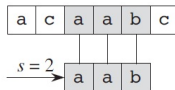
Ejemplo



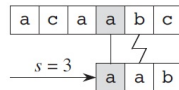
(a)



(b)



(c)



(d)

Algoritmo de Rabin Karp

Este algoritmo, utiliza la matemática teórica para comparar 2 números en módulo de un tercero. Para esto asume que cada carácter del texto es un número entre 0 y 9.

Notación:

- Texto $\rightarrow T[1..n] \rightarrow$ Valor decimal del texto del mismo tamaño que el patrón t_s .

Algoritmo de Rabin Karp

Este algoritmo, utiliza la matemática teórica para comparar 2 números en módulo de un tercero. Para esto asume que cada carácter del texto es un número entre 0 y 9.

Notación:

- Texto $\rightarrow T[1..n] \rightarrow$ Valor decimal del texto del mismo tamaño que el patrón t_s .
- Patrón $\rightarrow P[1..m] \rightarrow$ Valor decimal del patrón es p .

Algoritmo de Rabin Karp

Este algoritmo, utiliza la matemática teórica para comparar 2 números en módulo de un tercero. Para esto asume que cada carácter del texto es un número entre 0 y 9.

Notación:

- Texto $\rightarrow T[1..n] \rightarrow$ Valor decimal del texto del mismo tamaño que el patrón t_s .
- Patrón $\rightarrow P[1..m] \rightarrow$ Valor decimal del patrón es p .
- $d \rightarrow$ cantidad de números en el alfabeto.
En general $[0, 1, \dots, d-1]$

Algoritmo de Rabin Karp

Este algoritmo, utiliza la matemática teórica para comparar 2 números en módulo de un tercero. Para esto asume que cada carácter del texto es un número entre 0 y 9.

Notación:

- Texto $\rightarrow T[1..n] \rightarrow$ Valor decimal del texto del mismo tamaño que el patrón t_s .
- Patrón $\rightarrow P[1..m] \rightarrow$ Valor decimal del patrón es p .
- $d \rightarrow$ cantidad de números en el alfabeto.
En general $[0, 1, \dots, d-1]$
- $q \rightarrow$ Número primo con el que se calcula el módulo para hacer las comparaciones.

Preprocesamiento

Cálculo de p en $O(m)$ usando la regla de Horner

$$p = P[m] + 10(P[m-1] + 10(P[m-2] + \dots 10(P[2] + 10(P[2] + 10P[1])\dots))$$

valor de t_0 puede ser calculado desde $T[1..m]$ en $O(m)$. Para calcular t_1, t_2, \dots, t_{n-m} en tiempo $O(n-m)$ es suficiente observar que t_{s+1} puede ser calculado desde t_s en un tiempo constante:

Preprocesamiento

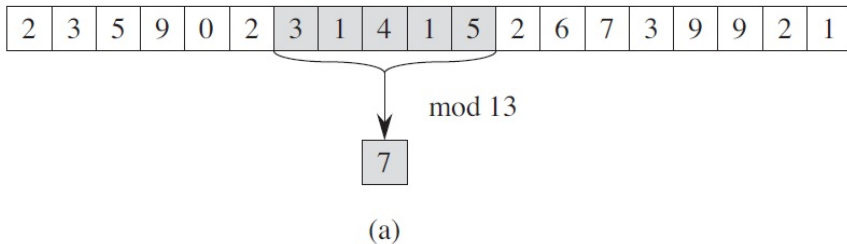
Cálculo de p en $O(m)$ usando la regla de Horner

$$p = P[m] + 10(P[m-1] + 10(P[m-2] + \dots 10(P[2] + 10(P[2] + 10P[1])\dots))$$

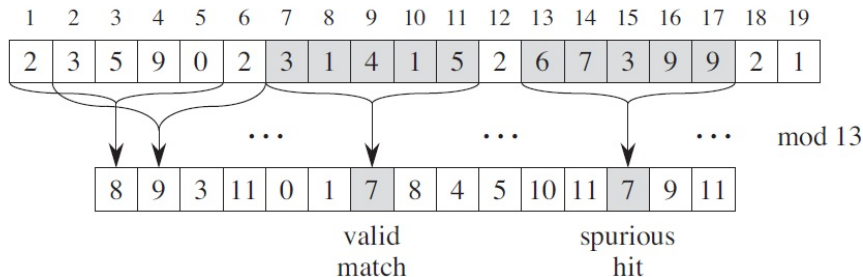
valor de t_0 puede ser calculado desde $T[1..m]$ en $O(m)$. Para calcular t_1, t_2, \dots, t_{n-m} en tiempo $O(n-m)$ es suficiente observar que t_{s+1} puede ser calculado desde t_s en un tiempo constante:

$$t_{s+1} = 10(t_s - 10^{m-1}T[s+1]) + T[s+m+1]$$

Ejemplo - Parte 1/3

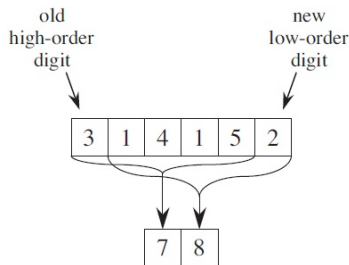


Ejemplo - Parte 2/3



(b)

Ejemplo - Parte 3/4



(c)

$$\begin{aligned}
 14152 &\equiv (31415 - 3 \cdot 10000) \cdot 10 + 2 \pmod{13} \\
 &\equiv (7 - 3 \cdot 3) \cdot 10 + 2 \pmod{13} \\
 &\equiv 8 \pmod{13}
 \end{aligned}$$

Diagram illustrating the Rabin-Karp algorithm's rolling hash calculation. The sequence 14152 is shown. The first digit '3' is labeled 'old high-order digit'. The last digit '2' is labeled 'new low-order digit'. A bracket under the digits 1, 4, 1, 5 indicates they are shifted one position to the left. Below this, the digits 7 and 8 are shown in a box, representing the new hash value.

Algoritmo del Automata Finito

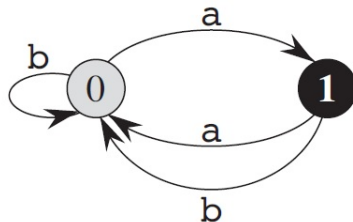
Este algoritmo se basa en la utilización de un autómata para encontrar un patrón. A medida que vamos leyendo caracteres, vamos avanzando de estado, y una vez que llegamos a un estado final significa que hubo matching.

Algoritmo del Automata Finito

Este algoritmo se basa en la utilización de un autómata para encontrar un patrón. A medida que vamos leyendo caracteres, vamos avanzando de estado, y una vez que llegamos a un estado final significa que hubo matching.

state	input	
	a	b
0	1	0
1	0	0

(a)



(b)

Notación de un Autómata Finito

Un autómata finito M se define mediante la tupla :

Notación de un Automata Finito

Un autómata finito M se define mediante la tupla :

$$\langle Q, q_0, A, \Sigma, \delta \rangle$$

donde :

Notación de un Automata Finito

Un autómata finito M se define mediante la tupla :

$$\langle Q, q_0, A, \Sigma, \delta \rangle$$

donde :

- Q es un conjunto finito de estados,

Notación de un Automata Finito

Un autómata finito M se define mediante la tupla :

$$\langle Q, q_0, A, \Sigma, \delta \rangle$$

donde :

- Q es un conjunto finito de estados,
- $q_0 \in Q$, es el estado inicial

Notación de un Automata Finito

Un autómata finito M se define mediante la tupla :

$$\langle Q, q_0, A, \Sigma, \delta \rangle$$

donde :

- Q es un conjunto finito de estados,
- $q_0 \in Q$, es el estado inicial
- $A \subseteq Q$, conjunto de estados finales

Notación de un Automata Finito

Un autómata finito M se define mediante la tupla :

$$\langle Q, q_0, A, \Sigma, \delta \rangle$$

donde :

- Q es un conjunto finito de estados,
- $q_0 \in Q$, es el estado inicial
- $A \subseteq Q$, conjunto de estados finales
- Σ , alfabeto de entrada

Notación de un Automata Finito

Un autómata finito M se define mediante la tupla :

$$\langle Q, q_0, A, \Sigma, \delta \rangle$$

donde :

- Q es un conjunto finito de estados,
- $q_0 \in Q$, es el estado inicial
- $A \subseteq Q$, conjunto de estados finales
- Σ , alfabeto de entrada
- $\delta : Q \times \Sigma \rightarrow Q$ es la función de transición

Notación de un Automata Finito

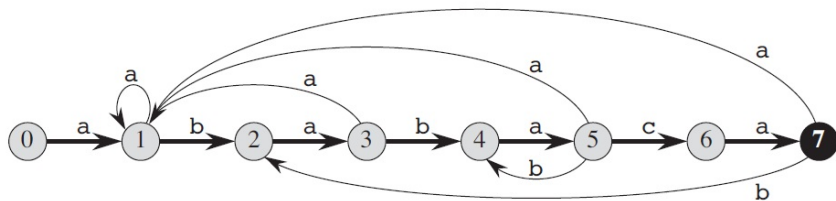
Un autómata finito M se define mediante la tupla :

$$\langle Q, q_0, A, \Sigma, \delta \rangle$$

donde :

- Q es un conjunto finito de estados,
- $q_0 \in Q$, es el estado inicial
- $A \subseteq Q$, conjunto de estados finales
- Σ , alfabeto de entrada
- $\delta : Q \times \Sigma \rightarrow Q$ es la función de transición
- $\phi(\varepsilon) = q_0$ $\phi(wa) = \delta(\phi(w), a)$ (transiciones sobre cadenas)

Ejemplo - Parte 1/2



(a)

Ejemplo - Parte 2/2

state	input			P
	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

(b)

i	—	1	2	3	4	5	6	7	8	9	10	11
$T[i]$	—	a	b	a	b	a	b	a	c	a	b	a
state $\phi(T_i)$	0	1	2	3	4	5	4	5	6	7	2	3

(c)

Algoritmo

FINITE-AUTOMATON-MATCHER(T, δ, m)

```
1   $n = T.length$ 
2   $q = 0$ 
3  for  $i = 1$  to  $n$ 
4       $q = \delta(q, T[i])$ 
5      if  $q == m$ 
6          print “Pattern occurs with shift”  $i - m$ 
```

Preprocesamiento

Función sufijo correspondiente a un patrón $P[1..m]$:

$$\sigma : \Sigma^* \rightarrow \{ 0,1,\dots,m \}$$

donde :

Preprocesamiento

Función sufijo correspondiente a un patrón $P[1..m]$:

$$\sigma : \Sigma^* \rightarrow \{ 0,1,\dots,m \}$$

donde :

- $\sigma(x)$ es la longitud del prefijo más largo de P que es sufijo de x :

Preprocesamiento

Función sufijo correspondiente a un patrón $P[1..m]$:

$$\sigma : \Sigma^* \rightarrow \{ 0,1,\dots,m \}$$

donde :

- $\sigma(x)$ es la longitud del prefijo más largo de P que es sufijo de x :
- $\sigma(x) = \{ k / P_k \sqsupseteq x \}$

Ejemplo

Si $P = ab$

- $\sigma(\varepsilon) =$

Ejemplo

Si $P = ab$

- $\sigma(\varepsilon) = 0$
- $\sigma(ccaca) =$

Ejemplo

Si $P = ab$

- $\sigma(\varepsilon) = 0$
- $\sigma(ccaca) = 1$
- $\sigma(ccaab) =$

Ejemplo

Si $P = ab$

- $\sigma(\varepsilon) = 0$
- $\sigma(ccaca) = 1$
- $\sigma(ccaab) = 2$

$P_0 = \varepsilon$ es sufijo de cualquier x

Dado un patrón P de longitud m

$$\sigma(x) = m$$

si y sólo si $P \sqsubset x$.

Definicion de AF

Definición del AF para $P[1..m]$

Definicion de AF

Definición del AF para $P[1..m]$

- Q es

Definicion de AF

Definición del AF para $P[1..m]$

- Q es $\{ 0, 1, \dots, m \}$

Definición de AF

Definición del AF para $P[1..m]$

- Q es $\{ 0, 1, \dots, m \}$
- q_0 es

Definición de AF

Definición del AF para $P[1..m]$

- Q es $\{ 0, 1, \dots, m \}$
- q_0 es 0

Definición de AF

Definición del AF para $P[1..m]$

- Q es $\{ 0, 1, \dots, m \}$
- q_0 es 0
- $A =$

Definición de AF

Definición del AF para $P[1..m]$

- Q es $\{ 0, 1, \dots, m \}$
- q_0 es 0
- $A = \{ m \}$

Definicion de AF

Definición del AF para $P[1..m]$

- Q es $\{ 0, 1, \dots, m \}$
- q_0 es 0
- $A = \{ m \}$
- $\delta(q, a) =$

Definición de AF

Definición del AF para $P[1..m]$

- Q es $\{ 0, 1, \dots, m \}$
- q_0 es 0
- $A = \{ m \}$
- $\delta(q, a) = \sigma(P_q a)$

Definición de AF

Definición del AF para $P[1..m]$

- Q es $\{ 0, 1, \dots, m \}$
- q_0 es 0
- $A = \{ m \}$
- $\delta(q, a) = \sigma(P_q a)$
- El AF mantiene invariante las transiciones:

$$\phi(T_i) = \sigma(T_i)$$

Algoritmo

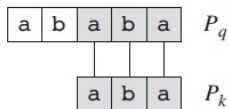
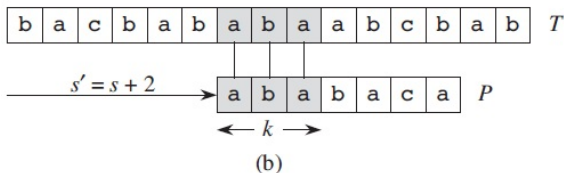
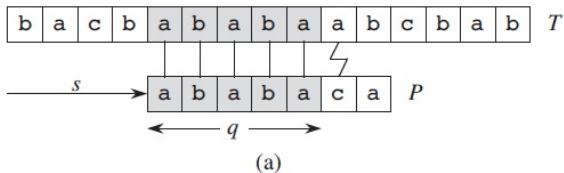
COMPUTE-TRANSITION-FUNCTION(P, Σ)

```
1   $m = P.length$ 
2  for  $q = 0$  to  $m$ 
3      for each character  $a \in \Sigma$ 
4           $k = \min(m + 1, q + 2)$ 
5          repeat
6               $k = k - 1$ 
7          until  $P_k \sqsupseteq P_q a$ 
8           $\delta(q, a) = k$ 
9  return  $\delta$ 
```

Algoritmo de Knuth Morris Pratt

Este algoritmo, utiliza un arreglo auxiliar $kmpNext[0 \dots m-1]$ el cual es generado con un costo de $O(m)$. Este arreglo brinda información adicional sobre el patrón a la hora de hacer las comparaciones, permite saber el corrimiento que hay que hacer en el patrón hasta la próxima comparación evitando preguntas innecesarias. Generando así un algoritmo eficiente de costo $O(n + m)$.

Ejemplo



Ejemplo

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

(a)

