

String Marching

Algoritmo do Knuth Morris Pratt

A.A. 2016-2017

Backtracking: string matching

Data una stringa P , chiamata pattern e una stringa (piu` lunga) T , chiamata testo, il problema del pattern matching esatto e` quello di trovare tutte le occorrenze del pattern P nel testo T .

Per esempio: $P = aba$, $T = bbabaxababay$
 P occorre in T tre volte, a partire dalle posizioni 3, 7 e 9.

A cosa serve?

- Programmi che trovano informazione testuale
- Programmi che cercano in cataloghi librari
- Internet browser
- Dizionari on-line
- Ricerche in basi dati che contengono parecchie centinaia di stringhe di DNA, RNA e aminoacidi.

Backtracking: string matching

Per esempio il genoma del batterio *Escherichia coli* di 4.6 milioni di nucleotidi ha milioni di repeat di lunghezza 12 e solo 8.000 di lunghezza massima 20 o più.
(il più frequente ACGCCGCATCCG di lunghezza 12 compare 94 volte).

La ricerca (usando GCG) in Genbank di una stringa di 30 caratteri ha richiesto 4 ore per verificare la sua non occorrenza. La ripetizione del test usando l'algoritmo di Boyer-Moore ha richiesto meno di 10 minuti, molti dei quali dedicati a swapping, meno di 1 minuto per la ricerca vera e propria.

Backtracking: string matching

Metodo naif: Allineare i caratteri di sinistra di P e T e confrontarli da sinistra a destra finché si trova un'occorrenza di P in T o si trovano due caratteri diversi. P viene spostato a destra di un carattere e il confronto ricomincia; termina quando il carattere più a sinistra di P supera quello più a sinistra di T.

P = aaa, T = aaaaaaaaaa vengono effettuati 24 confronti !

Sia $n = |P|$ e $m = |T|$, il numero di confronti nel caso peggiore è $\Theta(nm)$

P = AAAT, T = AAA...AAAAA se $|T| = 1$ milione la ricerca richiede circa 5 milioni di confronti e “disappointing” finisce senza successo.

La complessità può essere ridotta a $O(n + m)$

- Cercare di spostare P di più di un carattere quando si presenta un mismatch, mai però in modo da perdere un'occorrenza di P.
- Ridurre i confronti saltando alcune parti del pattern dopo lo shift.

Backtracking: string matching

```

0      1
1234567890123
T: xabxyabxyabxz
P: abxyabxz
*
abxyabxz
^^^^^^*
abxyabxz
*
abxyabxz
*
abxyabxz
*
abxyabxz
^^^^^^

```

20 confronti

```

0      1
1234567890123
T: xabxyabxyabxz
P: abxyabxz
*
abxyabxz
^^^^^^*
          abxyabxz
          ^^^^^^^

```

17 confronti

```

0      1
1234567890123
T: xabxyabxyabxz
P: abxyabxz
*
abxyabxz
^^^^^^*
          abxyabxz
          ^^^^^

```

14 confronti

Backtracking: string matching

E' facile disegnare un algoritmo di ricerca bruta.

Ricerca-bruta (P, T, n, m)

i, j, k: integer

$i \leftarrow 1$

$j \leftarrow 1$

$k \leftarrow 1$

while ($i \leq n$ and $j \leq m$)

if $T[i] = P[j]$

$i \leftarrow i+1; j \leftarrow j+1$

else $k \leftarrow k+1; i \leftarrow k; j \leftarrow 1$

if $j > m$ return k

else return i

Backtracking: string matching

Un backtrack meno brutale

$i = 14$
1 0 1 1 0 1 1 0 1 0 0 1 0 x x x x x x
1 0 0 1 0 1
 $j = 6$

$T[9] = P[1]$

$T[14] \neq P[6]$

$T[10] = P[2]$

$T[11] = P[3]$

$T[12] = P[4]$

e $P[1] P[2] = P[4] P[5]$

$T[13] = P[5]$

$T[12] T[13] = P[1] P[2]$



Si può far ripartire il confronto da $T[14]$ e $P[3]$

Essenziale: solo il pattern aiuta a calcolare l'indice da cui ricominciare il confronto, mentre sul testo non è necessario nessun backtracking.

Associamo ad ogni posizione j del pattern un indice che ricordi di quanto si deve “**saltare all'indietro**” se il confronto con il testo fallisce in quella posizione: $\text{back}[j]$.

Nell'esempio precedente abbiamo $\text{back}[6] = 3$

Backtracking: string matching

Algoritmo di Knuth - Morris - Pratt

n: lunghezza del testo, m:lunghezza del pattern

Calcola_back (P, back, m)

back[1] \leftarrow 0

k \leftarrow 0

for q \leftarrow 2 to m

while k > 0 and P[k+1] \neq P[q]

 k = back[k]

 if P[k+1] = P[q]

 k \leftarrow k + 1

 back[q] \leftarrow k

return

KMP (P, T, n, m)

Calcola_back (P, back, m)

 q \leftarrow 0

for i \leftarrow 1 to n

while q > 0 and P[q+1] \neq T[i]

 q \leftarrow back[q]

if P[q+1] = T[i]

 q \leftarrow q+1

if q = m

 print "pattern at i"

 q \leftarrow back[q]

return

Backtracking: string matching

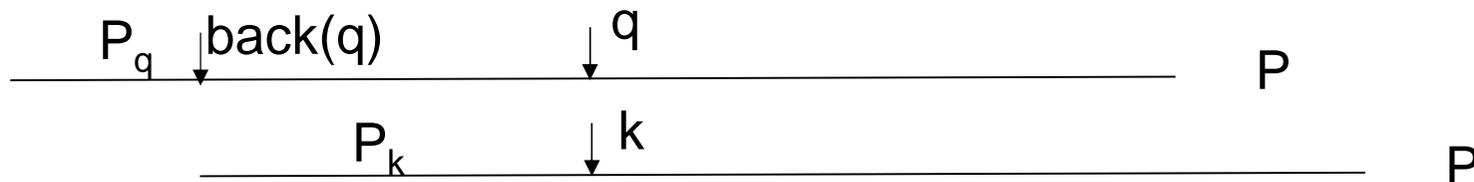
Perchè funziona? Proviamo a dare una spiegazione informale

Se P è una stringa chiamiamo P_q la sottostringa dei primi q caratteri
Indichiamo con $x \supseteq y$ se $x = uy$, cioè y è la parte terminale di x

Calcolando il vettore back su un pattern P si ha che:

$$\text{back}[q] = \max \{k \mid k < q \text{ and } P_q \supseteq P_k\}$$

In forma grafica:



Backtracking: string matching

Calcola_back (P, back, m)

back[1] ← 0

k ← 0

for q ← 2 to m

while k > 0 and P[k+1] ≠ P[q]

 k = back[k]

 if P[k+1] = P[q]

 k ← k + 1

 back[q] ← k

return 1

Abbreviamo *back* con *b*

Quando il controllo passa dal punto 1 si ha sempre:

$\forall q' \leq q :$

$b[q'] = \max \{k \mid k < q' \text{ and } P_{q'} \supseteq P_k\}$

Backtracking: string matching

Esempio

pattern = **a a c a b a a c c**
 1 2 3 4 5 6 7 8 9

back[1] = 0 back[7] = 2

back[2] = 1 back[8] = 3

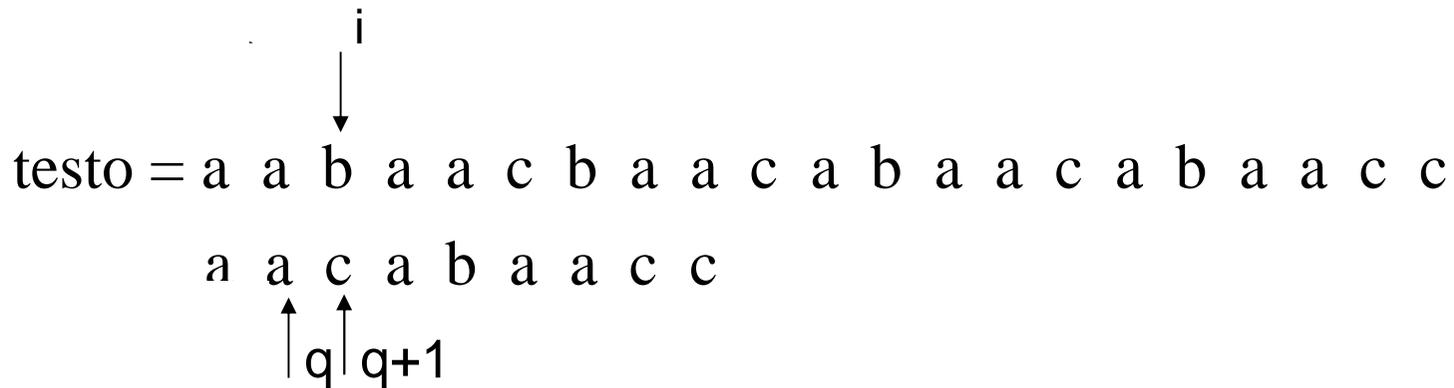
back[3] = 0 back[9] = 0

back[4] = 1

back[5] = 0

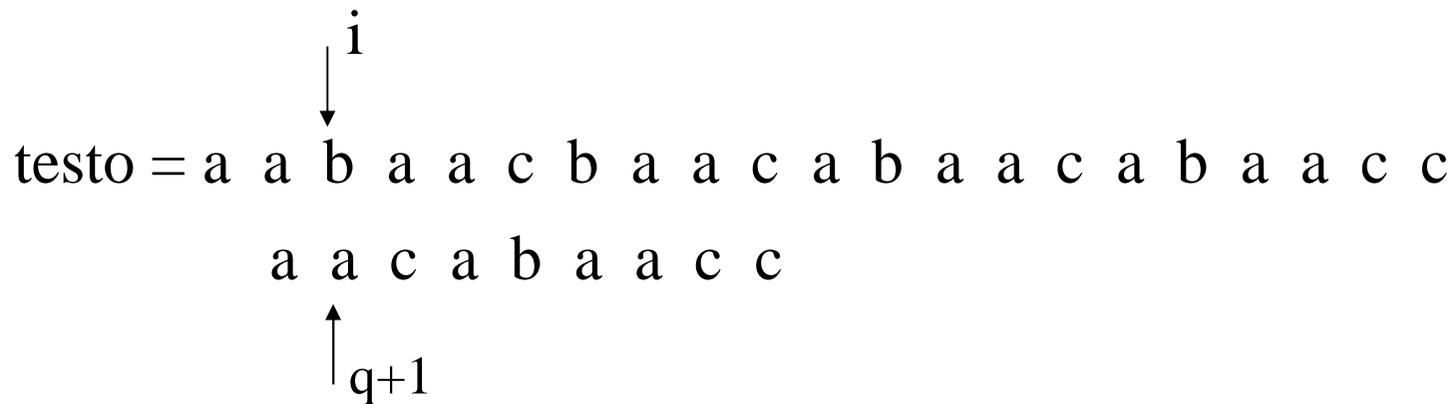
back[6] = 1

Backtracking: string matching



- back[1] = 0
- back[2] = 1
- back[3] = 0
- back[4] = 1
- back[5] = 0
- back[6] = 1
- back[7] = 2
- back[8] = 3
- back[9] = 0

back[2] = 1 “il confronto deve avvenire tra il primo carattere del pattern e il carattere successivo del testo, cioè tra gli elementi in posizione i e $q = 1$ ”



back[1] = 0 “il confronto fallisce e si resta nel while. si deve ricominciare dal primo carattere del pattern” e dal carattere successivo dell'input

Backtracking: string matching

$i \downarrow$
testo = a a b a a c b a a c a b a a c a b a a c c
a a c a b a a c c
 $\uparrow q+1$

back[1] = 0

back[2] = 1

back[3] = 0

back[4] = 1

back[5] = 0

back[6] = 1

back[7] = 2

back[8] = 3

back[9] = 0

back[3] = 0 “il confronto deve avvenire tra il primo carattere del pattern e il carattere successivo del testo, cioè tra gli elementi in posizione $i+1$ e 1”

$i \downarrow$
testo = a a b a a c b a a c a b a a c a b a a c c
a a c a b a a c c
 $\uparrow q+1$

back[8] = 3 “il confronto deve ricominciare da i con $q = 3$ nel pattern”

Backtracking: string matching

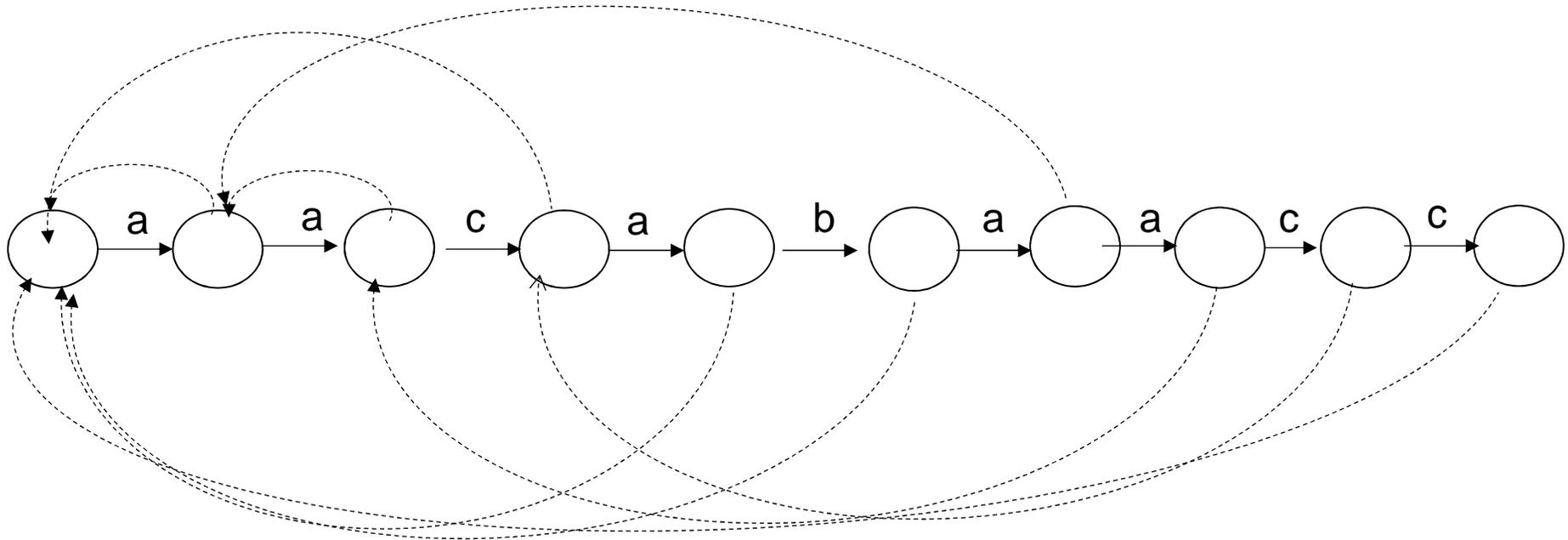
testo = a a b a a c b a a c a b a a c a b a a c c
a a c a b a a c c

↓ i
↑ q+1

il pattern è trovato e si finisce

Backtracking: string matching

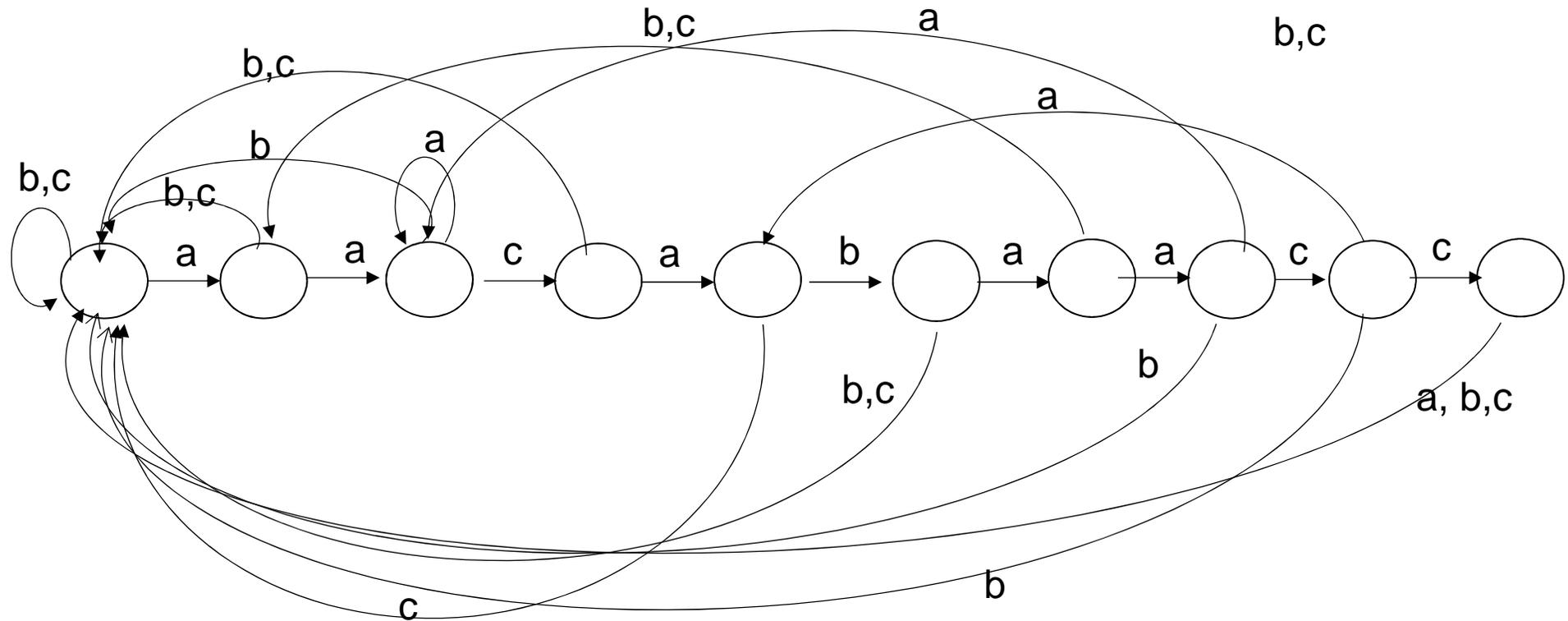
si può vedere come una specie di automa.



le transizioni all'indietro vengono eseguite per tutti quei simboli che non consentono di andare avanti, e si riparte (con lo stesso simbolo) dallo stato di arrivo.

Backtracking: string matching

Si può migliorare trasformandolo in un vero automa.



le transizioni all'indietro vengono eseguite per tutti quei simboli che non consentono di andare avanti, e si riparte (con lo stesso simbolo) dallo stato di arrivo.

Backtracking: string matching

Un altro esempio.

internal representation

j	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3	1	5	1
	B	0	2	0	4	4
	C	0	0	0	0	6

mismatch
transition
(back up)

match
transition
(increment)

graphical representation

