

Lezione 20 15/12/2023

Q-intervallo

$Q \rightarrow$ stringa definita su Σ

DEFINIZIONE1 (rispetto alla BWT)

Intervallo $[b,e)$ di posizioni della BWT che contengono i simboli che precedono i suffissi che condividono la stringa Q come prefisso

Sigma qui è l'alfabeto originale del testo, quindi senza il dollaro.

b è la posizione di inizio ed e è la posizione di fine, dove e è compresa, mentre b è esclusa.

DEFINIZIONE2 (rispetto al Suffix Array)

Intervallo $[b,e)$ di posizioni sul SA che contengono gli indici dei suffissi che condividono la stringa Q come prefisso

Q-intervallo per $Q=\epsilon \rightarrow [1,n+1)$

Se la stringa Q è la stringa vuota, l'epsilon-intervallo è quello che comprende tutta la BWT.

Esempio

$\Sigma = \{\$, a, c, g, t\}$

$T = acaaacatat\$$

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

$[2,8) \rightarrow$ a-intervallo

$[9,10) \rightarrow$ cat-intervallo

$[4,6) \rightarrow$ aca-intervallo = ac-intervallo

$[1,12) \rightarrow \epsilon$ -intervallo

Esempio

In un Q-intervallo $[b,e]$:

- ✓ $S[b], S[b+1], \dots, S[e-1]$
→ indici dei suffissi che hanno Q come prefisso
- ✓ $B[b], B[b+1], \dots, B[e-1]$
→ simboli che precedono i suffissi che hanno Q come prefisso

In un Q-intervallo $[b,e]$:

- ✓ $S[b], S[b+1], \dots, S[e-1]$
→ **occorrenze esatte di Q nel testo T**
- ✓ $B[b], B[b+1], \dots, B[e-1]$
→ **simboli che precedono le occorrenze esatte di Q nel testo T**

Quindi un q-intervallo rappresenta sempre le posizioni delle occorrenze esatte di Q nel testo.

Numero di occorrenze di Q in T → $(e-b)$

Nell'intervallo $[b,e]$

Esempio

L

$\Sigma = \{\$, a, c, g, t\}$

	S	B
1	11	t
2	3	c
3	4	a
4	1	\$
5	5	a
6	9	t
7	7	c
8	2	a
9	6	a
10	10	a
11	8	a

T = acaaacatat\$

↑ ↑

$[4,6) \rightarrow$ **aca**-intervallo

Le occorrenze di **aca** in T:

- sono in numero pari a $6-4=2$
- iniziano nelle posizioni: 1, 5
- sono precedute dai simboli: \$, a


Backward extension

La *backward extension* di un Q-intervallo $[b,e]$ con un simbolo σ è il σ Q-intervallo, cioè l'intervallo relativo alla stringa ottenuta concatenando il simbolo σ con Q

Prende in input un Q-intervallo e un simbolo, il risultato dell'operazione è il sigmaQ-intervallo, estendendo a sinistra la stringa aggiungendo sigma.


Esempio

	S	B
1	11 \$	t
2	3 <u>a</u> aacatat\$	c
3	4 <u>a</u> acatat\$	a
4	1 <u>a</u> caaacatat\$	\$
5	5 <u>a</u> catat\$	a
6	9 <u>a</u> t\$	t
7	7 <u>a</u> atat\$	c
8	2 <u>c</u> aaacatat\$	a
9	6 <u>c</u> atat\$	a
10	10 <u>t</u> \$	a
11	8 <u>t</u> at\$	a

[2,8) → a-intervallo
 backward extension con **c**
 [8,10) → **c**a-intervallo

	S	B
1	11 \$	t
2	3 <u>a</u> aacatat\$	c
3	4 <u>a</u> acatat\$	a
4	1 <u>a</u> caaacatat\$	\$
5	5 <u>a</u> catat\$	a
6	9 <u>a</u> t\$	t
7	7 <u>a</u> atat\$	c
8	2 <u>c</u> aaacatat\$	a
9	6 <u>c</u> atat\$	a
10	10 <u>t</u> \$	a
11	8 <u>t</u> at\$	a

	S	B
1	11 \$	t
2	3 <u>a</u> aacatat\$	c
3	4 <u>a</u> acatat\$	a
4	1 <u>a</u> caaacatat\$	\$
5	5 <u>a</u> catat\$	a
6	9 <u>a</u> t\$	t
7	7 <u>a</u> atat\$	c
8	2 <u>c</u> aaacatat\$	a
9	6 <u>c</u> atat\$	a
10	10 <u>t</u> \$	a
11	8 <u>t</u> at\$	a

[2,8) → a-intervallo
 backward extension con **g**
 intervallo vuoto (non esiste il **g**a-intervallo)

Ricerca esatta con BWT

Ricerca esatta di un pattern P lungo m

1. Si parte dal suffisso vuoto ε del pattern P
2. Si considera il Q-intervallo per $Q=\varepsilon$, cioè $[1, n+1)$, e si inizializza un indice di posizione $i=m$
3. Si effettua una *backward extension* con il simbolo $P[i]$ per ottenere il $P[i]$ Q-intervallo $[b_p, e_p)$ che fornisce le occorrenze del suffisso $P[i, m]$:
 - a) se il risultato è l'intervallo vuoto, allora P non ha occorrenze in T e la ricerca si ferma.
 - b) se il risultato non è l'intervallo vuoto, allora si decrementa di uno la posizione i e si ripete il punto 2 per l'intervallo ottenuto

2. m è l'ultimo simbolo del pattern

Alla prima iterazione fornirà le occorrenze dell'ultimo simbolo del pattern.

Devo distinguere il risultato che ottengo, se il risultato è vuoto allora mi fermo, perché non ci sono occorrenze.

Altrimenti si continua ad estendere all'indietro l'intervallo, fino a quando si raggiunge il primo simbolo del pattern.

- b) se il risultato non è l'intervallo vuoto
- se i è maggiore di 1, allora si decrementa di uno la posizione i e si ripete il punto 2 per l'intervallo ottenuto
 - se i è uguale a 1, significa che è stato trovato il Q-intervallo $[b_p, e_p]$ per $Q=P$ e vengono prodotte in output le occorrenze esatte $S[b_p], S[b_p + 1], \dots, S[e_p - 1]$

se $i > 1$, continuo ad estendere.

se $i = 1$ vuol dire che sono arrivato al primo simbolo del pattern con un intervallo non vuoto, quindi ho trovato le occorrenze del pattern nel testo.

Se P occorre in $T \rightarrow$ esattamente m iterazioni

✓ **Prima iterazione:** *backward extension* dell' ε -intervallo con il simbolo $P[m]$

$\rightarrow P[m]$ -intervallo

✓ **Iterazione intermedia:** *backward extension* del $P[i+1, m]$ -intervallo con il simbolo $P[i]$

$\rightarrow P[i, m]$ -intervallo

✓ **Ultima iterazione (m-esima):** *backward extension* del $P[2, m]$ -intervallo con il simbolo $P[1]$

$\rightarrow P$ -intervallo (non vuoto)

Se P non occorre in $T \rightarrow$ numero di iterazioni $\leq m$

✓ **Ultima iterazione:** *backward extension* del $P[p+1, m]$ -intervallo con il simbolo $P[p]$

$\rightarrow P[p, m]$ -intervallo vuoto

Massimo suffisso che occorre in $T \rightarrow P[p+1, m]$

Numero di iterazioni = $|P[p+1, m]| + 1$

Esempio

$T = \text{acaaacatat\$}$
 $P = \text{acat}$

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

$P = \text{acat}$

$[b, e] \leftarrow \varepsilon$ -intervallo $[1, 12]$

Iterazione 1

backward extension di $[1, 12]$ con $P[4]=t$
 $\rightarrow t$ -intervallo

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

P = acat

[b,e) ← ε-intervallo [1,12)

Iterazione 1

backward extension di [1,12) con P[4]=t
→ t-intervallo → [10,12)

Iterazione 2

backward extension di [10,12) con P[3]=a
→ at-intervallo → [6,8)

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

P = acat

[b,e) ← ε-intervallo [1,12)

Iterazione 1

backward extension di [1,12) con P[4]=t
→ t-intervallo → [10,12)

Iterazione 2

backward extension di [10,12) con P[3]=a
→ at-intervallo → [6,8)

Iterazione 3

backward extension di [6,8) con P[2]=c
→ cat-intervallo → [9,10)

	S	B
1	11 \$	t
2	3 aaacatat\$	c
3	4 aacatat\$	a
4	1 acaaacatat	\$
5	5 acatat\$	a
6	9 at\$	t
7	7 atat\$	c
8	2 caaacatat\$	a
9	6 catat\$	a
10	10 t\$	a
11	8 tat\$	a

P = acat

[b,e) ← ε-intervallo [1,12)

Iterazione 1

backward extension di [1,12) con P[4]=t
→ t-intervallo → [10,12)

Iterazione 2

backward extension di [10,12) con P[3]=a
→ at-intervallo → [6,8)

Iterazione 3

backward extension di [6,8) con P[2]=c
→ cat-intervallo → [9,10)

Iterazione 4

backward extension di [9,10) con P[1]=a
→ acat-intervallo → [5,6)

	S	F	B
1	11 \$		t
2	3 aaacatat\$		c
3	4 aacatat\$		a
4	1 acaaacatat		\$
5	5 acatat\$		a
6	9 at\$		t
7	7 atat\$		c
8	2 caaacatat\$		a
9	6 catat\$		a
10	10 t\$		a
11	8 tat\$		a

Il P-intervallo è [5,6)

→ 5 è l'unica occorrenza di P in T

Esempio

	S	F	B
1	11	\$	t
2	3	aaacatat\$	c
3	4	aacatat\$	a
4	1	acaaacatat	\$
5	5	acatat\$	a
6	9	at\$	t
7	7	atat\$	c
8	2	caaacatat\$	a
9	6	catat\$	a
10	10	t\$	a
11	8	tat\$	a

P = a**g**at

$[b,e) \leftarrow \varepsilon$ -intervallo $[1,12)$

Iterazione 1

backward extension di $[1,12)$ con $P[4]=t$
 $\rightarrow t$ -intervallo $\rightarrow [10,12)$

Iterazione 2

backward extension di $[10,12)$ con $P[3]=a$
 $\rightarrow at$ -intervallo $\rightarrow [6,8)$

Iterazione 3

backward extension di $[6,8)$ con $P[2]=c$
 $\rightarrow cat$ -intervallo $\rightarrow [9,10)$

~~Iterazione 4~~

~~backward extension di $[9,10)$ con $P[1]=a$
 $\rightarrow acat$ -intervallo $\rightarrow [5,6)$~~

Calcolo della backward extension

QUESTIONE: come trovare, dato un Q-intervallo $[b,e)$ e un simbolo σ , il σQ -intervallo?

	S	F	B
1	11	\$	t
2	3	a	c
3	4	a	a
4	1	a	\$
5	5	a	a
6	9	a	t
7	7	a	c
8	2	c	a
9	6	c	a
10	10	t	a
11	8	t	a

QUESTIONE: come trovare, dato il Q-intervallo $[2,8)$ per $Q=a$ e il simbolo c , il cQ -intervallo?

$i_1=2 \rightarrow B[i_1]=c$ precede il 3-suffisso

$i_2=7 \rightarrow B[i_2]=c$ precede il 7-suffisso

il 3-suffisso ha Q come prefisso

il 7-suffisso ha Q come prefisso

$B[i_1] + 3$ -suffisso \rightarrow suffisso che ha cQ come prefisso

$B[i_2] + 7$ -suffisso \rightarrow suffisso che ha cQ come prefisso

$j_1 = LF(i_1=2) = 8$

$j_2 = LF(i_2=7) = 9$

7	7	a	c
8	2	c	a
9	6	c	a
10	10	t	a
11	8	t	a

$B[i_1] + 3$ -suffisso \rightarrow suffisso che ha cQ come prefisso

$B[i_2] + 7$ -suffisso \rightarrow suffisso che ha cQ come prefisso

$j_1 = LF(i_1=2) = 8$

$j_2 = LF(i_2=7) = 9$

$[8,10)$ è il cQ -intervallo

Q-intervallo $[b,e)$ e simbolo σ :

- ✓ siano $i_1 < i_2 < \dots < i_k$ le k posizioni in $[b,e)$ tali che:

$$B[i_p] = \sigma \quad 1 \leq p \leq k$$

- ✓ $B[i_p]$ precede $S[i_p]$ -suffisso per $1 \leq p \leq k$
- ✓ $S[i_p]$ -suffisso ha Q come prefisso per $1 \leq p \leq k$
- ✓ $B[i_p] + S[i_p]$ -suffisso per $1 \leq p \leq k$ è il suffisso che inizia con $B[i_p]$ e che ha σQ come prefisso
- ✓ $B[i_p] + S[i_p]$ -suffisso per LF-mapping avrà una posizione nel Suffix Array data da $j_p = LF(i_p)$
- ✓ $j_1, j_2, \dots, j_p, \dots, j_k \rightarrow [j_1, j_k+1)$ è il σQ -intervallo

Procedura **Backward_extend**(b, e, σ)

```
 $i_1 \leftarrow$  più piccola posizione  $\geq b$  tale che  $B[i_1] = \sigma$   
 $i_k \leftarrow$  più grande posizione  $< e$  tale che  $B[i_k] = \sigma$   
  
return  $[LF(i_1), LF(i_k)+1)$ 
```

Procedura **Search_pattern**(P, S)

```
 $n \leftarrow |S|$   
 $[b,e) \leftarrow [1, n+1)$   
 $i \leftarrow |P|$   
while  $[b,e)$  is not null and  $i \geq 1$  do  
   $\sigma \leftarrow P[i]$   
   $[b,e) \leftarrow$  Backward_extend( $b, e, \sigma$ )  
   $i \leftarrow i-1$   
if  $[b,e)$  is not null then  
  output  $S[b], S[b+1], \dots, S[e-1]$ 
```

Per il momento...

Complessità $\neq O(m)$

FM-index

è una rappresentazione della BWT in forma numerica.

FM-index = funzione **C** + funzione **Occ**

C: $\Sigma \rightarrow \{0, 1, 2, 3, \dots, n\}$

C(σ) = numero di simboli della BWT che sono $< \sigma$

La funzione **C** mappa i simboli dell'alfabeto (dollaro compreso) in un numero da 0 a n .

Occ: $\{1, 2, 3, \dots, n+1\} \times \Sigma \rightarrow \{0, 1, 2, 3, \dots, n\}$

Occ(i, σ) = numero di simboli σ in $B[1, i-1]$

Occ mappa ogni coppia che si può formare con un numero da 1 a $n+1$ e un simbolo, ad un intero da 0 a n .

Queste due funzioni rappresentano la BWT. Questo è un self index, quindi potremmo tenere queste due funzioni senza la BWT, perché si può ricostruire partendo da queste sia la BWT che il testo.

Costruiamo la funzione **C**:

	1	2	3	4	5	6	7	8	9	
										\$ < a < c < g < t
T	g	g	t	c	a	g	t	c	\$	

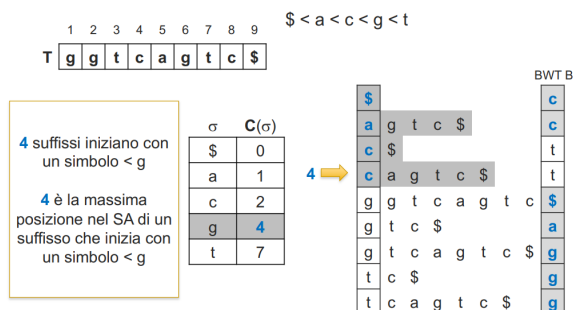
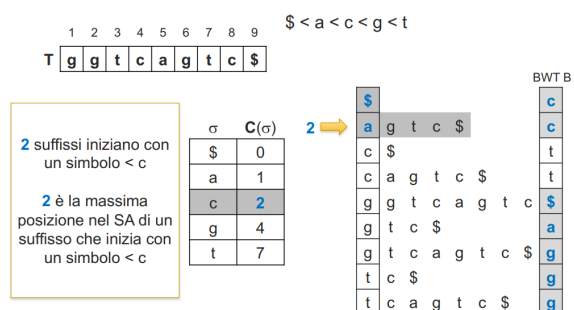
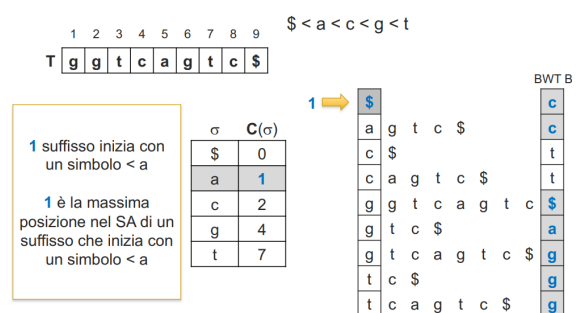
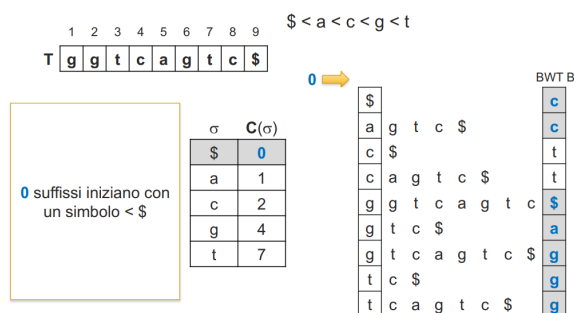
σ	$C(\sigma)$
\$	0
a	1
c	2
g	4
t	7

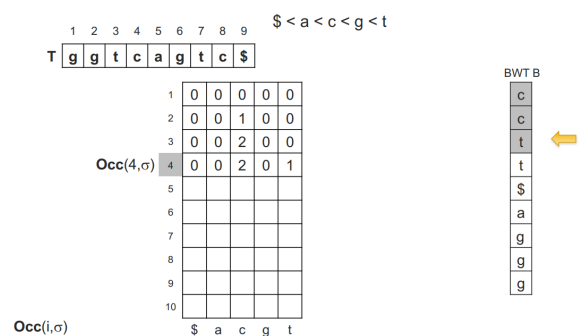
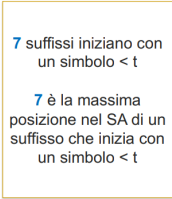
BWT B

c
c
t
t
\$
a
g
g
g

$C(\sigma)$ fornisce il numero di suffissi che iniziano con un simbolo lessicograficamente inferiore a σ .

Questo valore $C(\sigma)$ fornisce la massima posizione nel suffix array di un simbolo che inizia con un simbolo inferiore a σ .





	1	2	3	4	5	6	7	8	9		\$ < a < c < g < t
T	g	g	t	c	a	g	t	c	\$		

1	0	0	0	0	0
2	0	0	1	0	0
3	0	0	2	0	0
4	0	0	2	0	1
5	0	0	2	0	2
6	1	0	2	0	2
7	1	1	2	0	2
8	1	1	2	1	2
9	1	1	2	2	2
10	1	1	2	3	2

Occ(i,σ)

	\$	a	c	g	t
--	----	---	---	---	---

BWT B

c
c
t
t
\$
a
g
g
g



	1	1	2	1	2
8	1	1	2	1	2
9	1	1	2	2	2
10	1	1	2	3	2

\$ a c g t

Distribuzione dei simboli in B

L'ultima riga della tabella ci dice che abbiamo 1 dollaro, 1a, 2c, 3g e 2t. Di conseguenza potremmo ricavare la funzione C dal quest'ultima riga.

σ	C(σ)
\$	0
a	1
c	2
g	4
t	7

Procedura Build-FMindex(B)

$n \leftarrow |B|$

$C \leftarrow$ empty vector of size $|\Sigma|$

$Occ \leftarrow$ empty table $(n+1) \times |\Sigma|$

foreach σ **do**

$Occ[1, \sigma] \leftarrow 0$

for i from 1 to n **do**

foreach σ in Σ **do**

$Occ[i+1, \sigma] \leftarrow Occ[i, \sigma]$

$Occ[i+1, B[i]] \leftarrow Occ[i+1, B[i]] + 1$

$C[\$] \leftarrow 0$

//Scorro i simboli in ordine lessicografico

foreach σ in $\Sigma \setminus \{\$ \}$ **do**

$\sigma' \leftarrow$ simbolo immediatamente prima di σ in Σ

$C[\sigma] \leftarrow C[\sigma'] + Occ[n+1, \sigma']$

return (C, Occ)

Calcolo della LF-function

ha fatto qualche slides e esempio ma dovrebbe rifarlo la prossima volta