

# Lezione 21 20/12/2023

## LF-function

### Proprietà di LF-mapping

Il suffisso che inizia con l' $i$ -esimo simbolo  $\sigma$  della BWT  $B$  è l' $i$ -esimo suffisso che inizia con il simbolo  $\sigma$  nel Suffix Array

### LF-function

Funzione che fornisce la posizione lessicografica  $j$  del suffisso che inizia con il simbolo in posizione  $i$  della BWT, cioè  $B[i]$

$$j = LF(i)$$

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

Per  $i=8$ ,  $B[i] = g$  è il simbolo iniziale del secondo suffisso che inizia con  $g$ , dal momento che  $B[i]$  è la seconda  $g$  in  $B$

$LF(i)$  deve restituire la posizione  $j=6$  del suffisso **gtc\$** nel Suffix Array

$$6 = 4 + 2$$

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

$$6 = C(B[8]) + Occ(8, B[8]) + 1$$

$$j = C(B[i]) + Occ(i, B[i]) + 1$$

↑  
LF-function

Ci sono 4 suffissi che iniziano con un simbolo inferiore di  $g$  (lo leggo dalla funzione  $B$ ) e sommo 2 perché questa è la seconda  $g$  della BWT. Questo 2 lo leggo dalla funzione  $Occ$  (numero di prefissi fino alla posizione 7) + 1 perché la  $Occ$  va fino ad una posizione prima.

## Calcolo della backward extension

La backward extension di un Q-intervallo che inizia nella posizione  $b$  e finisce nella posizione  $e-1$ , aggiunge il simbolo  $\sigma$  davanti e avrà una posizione di inizio e fine diversa.

Q-intervallo  $[b, e)$ ,  $\sigma \rightarrow \sigma$ Q-intervallo  $[b', e')$

$$b' = LF(i_1)$$

$$e' = LF(i_k) + 1$$

$i_1 \rightarrow$  più piccola posizione in  $[b, e)$  tale che  $B[i_1] = \sigma$

$i_k \rightarrow$  più grande posizione in  $[b, e)$  tale che  $B[i_k] = \sigma$

**LF-Function:**

$$j = LF(i) = C(B[i]) + \text{Occ}(i, B[i]) + 1$$

$$b' = C(B[i_1]) + \text{Occ}(i_1, B[i_1]) + 1$$

$$e' = C(B[i_k]) + \text{Occ}(i_k, B[i_k]) + 1 + 1$$

$$b' = C(\sigma) + \text{Occ}(i_1, \sigma) + 1$$

$$e' = C(\sigma) + \text{Occ}(i_k, \sigma) + 1 + 1$$

In  $B[b, i_1-1]$  **non** esistono simboli uguali a  $\sigma$   
 $\Rightarrow$  il numero di simboli  $\sigma$  in  $B[1, i_1-1]$  è uguale al numero di simboli  $\sigma$  in  $B[1, b-1] \Rightarrow \text{Occ}(i_1, \sigma) = \text{Occ}(b, \sigma)$

perchè  $i_1$  è quella più piccola che contiene  $\sigma$ .

Quindi ci svincoliamo da  $i_1$ , lo sostituiamo con  $b$ .

facciamo un ragionamento simile sul secondo

1.  $\text{Occ}(i_k, \sigma) =$  numero di simboli  $\sigma$  in  $B[1, i_k-1]$   
2.  $B[i_k] = \sigma$   
 $\Rightarrow \text{Occ}(i_k, \sigma) + 1 = \text{Occ}(i_k+1, \sigma)$

In  $B[i_k+1, e-1]$  non esistono simboli uguali a  $\sigma$   
 $\Rightarrow$  il numero di simboli  $\sigma$  in  $B[1, i_k]$  è uguale al numero di simboli  $\sigma$  in  $B[1, e-1] \Rightarrow \text{Occ}(i_k+1, \sigma) = \text{Occ}(e, \sigma)$

~~$i_1 \rightarrow$  più piccola posizione in  $[b, e)$  tale che  $B[i_1] = \sigma$~~   
 ~~$i_k \rightarrow$  più grande posizione in  $[b, e)$  tale che  $B[i_k] = \sigma$~~

$$b' = C(\sigma) + \text{Occ}(b, \sigma) + 1$$

$$e' = C(\sigma) + \text{Occ}(e, \sigma) + 1$$

$e' = b' \Rightarrow$  il  $\sigma$ Q-intervallo è vuoto

**Procedura Backward\_extend( $b, e, \sigma$ )**

$$b' = C(\sigma) + \text{Occ}(b, \sigma) + 1$$

$$e' = C(\sigma) + \text{Occ}(e, \sigma) + 1$$

**return**  $[b', e')$

**Tempo  $O(1)$**

**Procedura Search\_pattern( $P, S$ )**

**begin**

$n \leftarrow |S|$

$[b, e) \leftarrow [1, n+1)$

$i \leftarrow |P|$

**while**  $[b, e)$  **is not** null **and**  $i \geq 1$  **do**

$\sigma \leftarrow P[i]$

$[b, e) \leftarrow \text{Backward\_extend}(b, e, \sigma)$

$i \leftarrow i-1$

**if**  $[b, e)$  **is not** null **then**

output  $S[b], S[b+1], \dots, S[e-1]$

**end**

Complessità  $\rightarrow O(m)$

Esercizio FM-index è un self-index

Ricostruisci la BWT, che sarà lunga 1 in meno rispetto alla Occ.

**Occ(i,σ)**

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

σ	C(σ)
\$	0
a	1
c	2
g	4
t	7

**BWT B**

1	
2	
3	
4	
5	
6	
7	
8	
9	

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

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

**BWT B**

1	c
2	
3	
4	
5	
6	
7	
8	
9	

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

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

**BWT B**

1	c
2	c
3	
4	
5	
6	
7	
8	
9	

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

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

**BWT B**

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

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

Vediamo cos'è cambiato da ogni riga rispetto alla riga precedente.

Esercizio 1

Si consideri la seguente funzione C. Dire se il b-intervallo (per Q=b) è vuoto. Se non è vuoto, specificare le posizioni di inizio e fine.

σ	C(σ)
\$	0
a	1
b	3
c	7
d	7

Funzione C

$\sigma$	$C(\sigma)$
\$	0
a	1
<b>b</b>	<b>3</b>
<b>c</b>	<b>7</b>
d	7

$C(c) - C(b) > 0$  e quindi esistono suffissi che iniziano con il simbolo b.

→ b-intervallo **non** è vuoto

Funzione C

$\sigma$	$C(\sigma)$
\$	0
a	1
<b>b</b>	<b>3</b>
<b>c</b>	<b>7</b>
d	7

Numero di suffissi che iniziano con  $\sigma < b \rightarrow C(b) = 3$

Posizione nel Suffix Array del primo suffisso che inizia con b →  $C(b) + 1 = 4$

Numero di suffissi che iniziano con  $\sigma < c \rightarrow C(c) = 7$

Posizione nel Suffix Array dell'ultimo suffisso che inizia con b →  $C(c) = 7$

Funzione C

**b-intervallo → [4,8)**

## Esercizio 2

Data la Burrows-Wheeler Transform  $B = \text{accgt\$ac}$  di un testo  $T$ , si richiede di specificare l'FM-index supponendo  $\Sigma = \{a, c, g, t\}$ . Calcolare poi tramite FM-index la posizione  $j$  nel Suffix Array del suffisso che inizia con il terzo simbolo della BWT.

Funzione C?

B
a
c
c
g
t
\$
a
c

$\sigma$	$C(\sigma)$
\$	0

B
a
c
c
g
t
\$
a
c

$\sigma$	$C(\sigma)$
\$	0
a	1

B
a
c
c
g
t
\$
a
c

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3

B

a
c
c
g
t
\$
a
c

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3
g	6
t	

B

a
c
c
g
t
\$
a
c

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3
g	6
t	7

Funzione **Occ?**

B

a
c
c
g
t
\$
a
c

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3
g	6
t	7

1	0	0	0	0	0
2					
3					
4					
5					
6					
7					
8					
9					
	\$	a	c	g	t

B

a
c
c
g
t
\$
a
c

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3
g	6
t	7

1	0	0	0	0	0
2	0	1	0	0	0
3	0	1	1	0	0
4	0	1	2	0	0
5	0	1	2	1	0
6	0	1	2	1	1
7	1	1	2	1	1
8	1	2	2	1	1
9	1	2	3	1	1
	\$	a	c	g	t

Posizione j nel SA del suffisso che inizia B[3]?

3 →

B

a
c
c
g
t
\$
a
c

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3
g	6
t	7

1	0	0	0	0	0
2	0	1	0	0	0
3	0	1	1	0	0
4	0	1	2	0	0
5	0	1	2	1	0
6	0	1	2	1	1
7	1	1	2	1	1
8	1	2	2	1	1
9	1	2	3	1	1
	\$	a	c	g	t

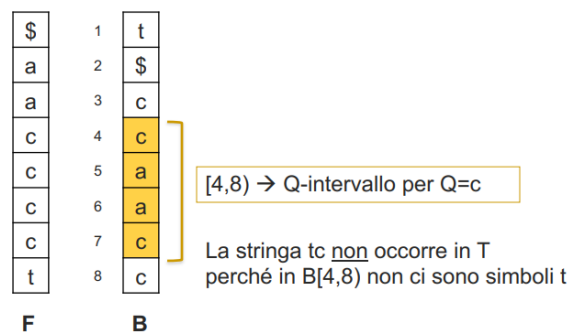
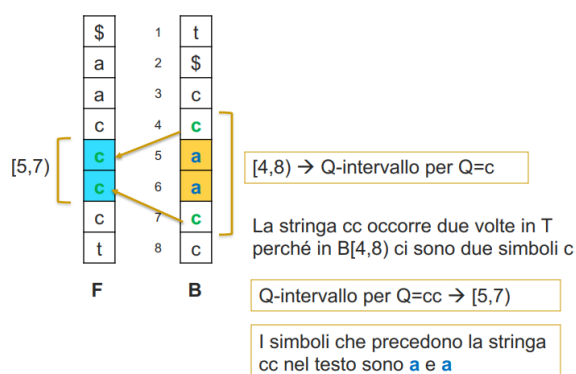
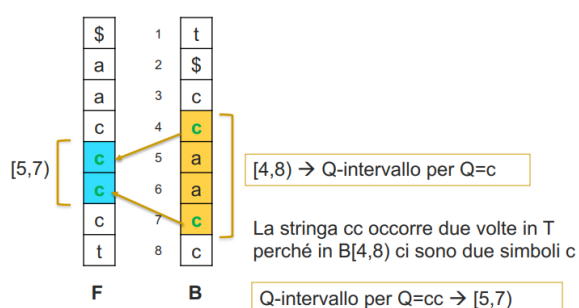
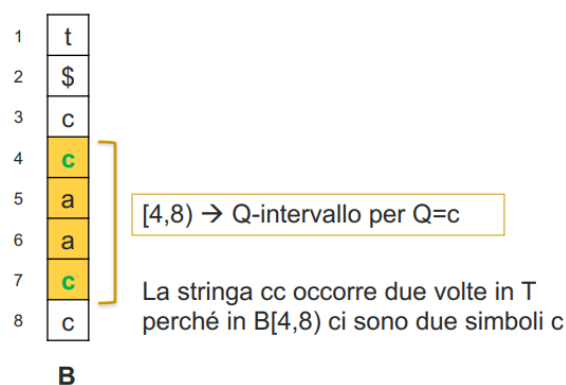
$$j = C(c) + \text{Occ}(3, c) + 1 = 3 + 1 + 1 = 5$$

## Esercizio 3

Data la BWT  $B = t\$ccaacc$  di un testo  $T$  e sapendo che il Q-intervallo  $[4,8)$  rappresenta la stringa  $Q=c$ , specificare utilizzando unicamente  $B$ :

- 1) quante volte la stringa  $cc$  occorre in  $T$ , e specificare il Q-intervallo per  $Q=cc$
- 2) quante volte la stringa  $tc$  occorre in  $T$ , e specificare il Q-intervallo per  $Q=tc$

Indicare inoltre quali sono i simboli che nel testo precedono le occorrenze di  $cc$  e le occorrenze di  $tc$ .



## Esercizio 4

Per un testo \$-terminato di lunghezza  $n=10$  la funzione **C** di FM-index è:

$$C(\$) = 0; C(a) = 1; C(c) = 4; C(g) = 6; C(t) = 6$$

1) Dire se una funzione **Occ** che per  $i = 11$  vale:

$$\mathbf{Occ}(11, \$) = 1; \mathbf{Occ}(11, a) = 2; \mathbf{Occ}(11, c) = 3; \mathbf{Occ}(11, g) = 2; \mathbf{Occ}(11, t) = 2$$

è compatibile con la funzione **C**. In caso contrario, specificare tutte le incongruenze.

2) Dire inoltre (in base alla funzione **C**) quanti sono i simboli **g** e quanti sono i simboli **t** nel testo.

**RISPOSTA1:** NO, perché

**Distribuzione dei simboli derivata dalla funzione C:**

$$C(a) - C(\$) = 1 \text{ simbolo } \$$$

$$C(c) - C(a) = 3 \text{ simboli } a$$

$$C(g) - C(c) = 2 \text{ simboli } c$$

$$C(t) - C(g) = 0 \text{ simboli } g$$

$$10 - C(t) = 4 \text{ simboli } t$$

**RISPOSTA2:**

$$\text{Numero di simboli } g \rightarrow 0$$

$$\text{Numero di simboli } t \rightarrow 4$$

che è diversa da quella specificata dalla funzione **Occ**

**RISPOSTA1:** NO, perché

$$a) \mathbf{Occ}(11, \$) + \mathbf{Occ}(11, a) = 3 \text{ è diverso da } C(c) = 4$$

$$b) \mathbf{Occ}(11, \$) + \mathbf{Occ}(11, a) + \mathbf{Occ}(11, c) + \mathbf{Occ}(11, g) = 8 \text{ è diverso da } C(t) = 6$$

**RISPOSTA2:**

$$\text{Numero di simboli } g \rightarrow C(t) - C(g) = 0$$

$$\text{Numero di simboli } t \rightarrow n - C(t) = 10 - 6 = 4$$



Esercizio 5

1	0	0	0	0	0
2	0	0	0	0	1
3	0	0	0	1	1
4	0	1	0	1	1
5	0	1	0	2	1
6	0	1	0	2	2
7	0	2	0	2	2
8	0	2	1	2	2
9	1	2	1	2	2
10	1	2	1	3	2
11	1	2	1	4	2
	\$	a	c	g	t

Data la funzione **Occ** di FM-index di un testo T, derivare:

- la funzione **C**
- il simbolo della BWT in posizione 4

7	0	2	0	2	2
8	0	2	1	2	2
9	1	2	1	2	2
10	1	2	1	3	2
11	1	2	1	4	2
	\$	a	c	g	t

$\sigma$	<b>C</b> ( $\sigma$ )
\$	0
a	1
c	
g	
t	

Funzione **C**

7	0	2	0	2	2
8	0	2	1	2	2
9	1	2	1	2	2
10	1	2	1	3	2
11	1	2	1	4	2
	\$	a	c	g	t

$\sigma$	<b>C</b> ( $\sigma$ )
\$	0
a	1
c	3
g	
t	

Funzione **C**

7	0	2	0	2	2
8	0	2	1	2	2
9	1	2	1	2	2
10	1	2	1	3	2
11	1	2	1	4	2
	\$	a	c	g	t

$\sigma$	<b>C</b> ( $\sigma$ )
\$	0
a	1
c	3
g	4
t	

Funzione **C**

7	0	2	0	2	2
8	0	2	1	2	2
9	1	2	1	2	2
10	1	2	1	3	2
11	1	2	1	4	2
	\$	a	c	g	t

$\sigma$	<b>C</b> ( $\sigma$ )
\$	0
a	1
c	3
g	4
t	8

Funzione **C**

1	0	0	0	0	0
2	0	0	0	0	1
3	0	0	0	1	1
4	0	1	0	1	1
5	0	1	0	2	1
6	0	1	0	2	2
7	0	2	0	2	2
8	0	2	1	2	2
9	1	2	1	2	2
10	1	2	1	3	2
11	1	2	1	4	2
	\$	a	c	g	t

Data la funzione **Occ** di FM-index di un testo T, derivare:

- la funzione **C**

- il simbolo della BWT in posizione 4

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3
g	4
t	8

$B[4] = g$

Funzione **C**

Inoltre, sapendo che  $[2,4]$  è il Q-intervallo per  $Q=a$ , dire sulla base di FM-index se la stringa aa occorre nel testo T.

La stringa aa occorre in T se in  $B[2,4]$  esiste almeno un simbolo a

In posizione 2 esiste un simbolo a?

1	0	0	0	0	0
2	0	0	0	0	1
3	0	0	0	1	1
4	0	1	0	1	1
5	0	1	0	2	1
6	0	1	0	2	2
7	0	2	0	2	2
8	0	2	1	2	2
9	1	2	1	2	2
10	1	2	1	3	2
11	1	2	1	4	2
	\$	a	c	g	t

Inoltre, sapendo che  $[2,4]$  è il Q-intervallo per  $Q=a$ , dire sulla base di FM-index se la stringa aa occorre nel testo T.

La stringa aa occorre in T se in  $B[2,4]$  esiste almeno un simbolo a

In posizione 2 esiste un simbolo a?  
**NO**,  $B[2]$  è uguale a g

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3
g	4
t	8

Funzione **C**

1	0	0	0	0	0
2	0	0	0	0	1
3	0	0	0	1	1
4	0	1	0	1	1
5	0	1	0	2	1
6	0	1	0	2	2
7	0	2	0	2	2
8	0	2	1	2	2
9	1	2	1	2	2
10	1	2	1	3	2
11	1	2	1	4	2
	\$	a	c	g	t

Inoltre, sapendo che  $[2,4]$  è il Q-intervallo per  $Q=a$ , dire sulla base di FM-index se la stringa aa occorre nel testo T.

La stringa aa occorre in T se in  $B[2,4]$  esiste almeno un simbolo a

In posizione 3 esiste un simbolo a?  
**YES**,  $B[3]$  è uguale ad a

La stringa aa occorre in T

$\sigma$	$C(\sigma)$
\$	0
a	1
c	3
g	4
t	8

Funzione **C**

## Esercizio 6

Si consideri la seguente funzione **Occ** di FM-index e si determini il testo originale

1	0	0	0	0	0
2	0	0	0	0	1
3	1	0	0	0	1
4	1	0	1	0	1
5	1	0	2	0	1
6	1	1	2	0	1
7	1	2	2	0	1
8	1	2	3	0	1
9	1	2	4	0	1
\$ a c g t					

BWT B?

1	0	0	0	0	0
2	0	0	0	0	1
3	1	0	0	0	1
4	1	0	1	0	1
5	1	0	2	0	1
6	1	1	2	0	1
7	1	2	2	0	1
8	1	2	3	0	1
9	1	2	4	0	1
\$ a c g t					

t

B

1	0	0	0	0	0
2	0	0	0	0	1
3	1	0	0	0	1
4	1	0	1	0	1
5	1	0	2	0	1
6	1	1	2	0	1
7	1	2	2	0	1
8	1	2	3	0	1
9	1	2	4	0	1
\$ a c g t					

t
\$

B

1	0	0	0	0	0
2	0	0	0	0	1
3	1	0	0	0	1
4	1	0	1	0	1
5	1	0	2	0	1
6	1	1	2	0	1
7	1	2	2	0	1
8	1	2	3	0	1
9	1	2	4	0	1
\$ a c g t					

t
\$
c

B

Ricostruzione di T

1	0	0	0	0	0
2	0	0	0	0	1
3	1	0	0	0	1
4	1	0	1	0	1
5	1	0	2	0	1
6	1	1	2	0	1
7	1	2	2	0	1
8	1	2	3	0	1
9	1	2	4	0	1
\$ a c g t					

\$
a
a
c
c
c
c
t

F

t
\$
c
c
a
a
c
c

B

T = accacct\$