

# LA PIU' LUNGA SOTTOSEQUENZA COMUNE (LONGEST COMMON SUBSEQUENCE : LCS)

## • APPLICAZIONI BIOLOGICHE

- CONFRONTO **DNA** DI DIVERSI ORGANISMI

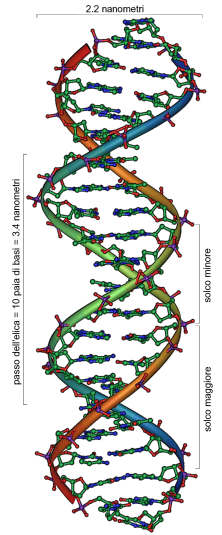
- **STRUTTURA DNA** : FORMATA DA UNA STRINGA DI MOLECOLE  
DETTE **BASI** (ca.  $3.2 \times 10^9$ )

- **BASI** :

- **A**DENINA (A)
- **C**ITOSINA (C)
- **G**UANINA (G)
- **T**IMINA (T)

- QUINDI IL DNA DI UN ORGANISMO PUO' ESSERE  
RAPPRESENTATO COME UNA STRINGA NELL'ALFABETO **{A, C, G, T}**

- LA **CORRELAZIONE** TRA DUE ORGANISMI PUO' ESSERE  
"MISURATA" CON IL GRADO DI **SOMIGLIANZA** DEI LORO DNA



Es.  $S_1 = \text{ACCGGTCGAGTGCGCGGAAGCCGGCCGAA}$   
 $S_2 = \text{GTCGTTCGGAATGCCGTTGCTCTGTAAA}$   
 $S_3 = \text{GTCGTCGGAAGCCGGCCGAA}$

VI SONO VARI CRITERI DI SOMIGLIANZA :

- UNO DEI DUE FILAMENTI E' SOTTOSTRINGA DELL'ALTRO  
[STRING MATCHING]
- IL NUMERO DI MODIFICHE DA FARE AD UNO DEI  
DUE FILAMENTI DI DNA PER OTTENERE L'ALTRO  
E' "PICCOLO" [DISTANZA DI EDIT]
- C'E' UN "LUNGO" FILAMENTO,  $S_3$ , LE CUI BASI  
COMPATONO SIA IN  $S_1$ , CHE IN  $S_2$  (NELLO  
STESSO ORDINE, ANCHE SE IN POSIZIONI NON  
NECESSARIAMENTE CONSECUTIVE) [LCS]

DEF UNA SOTTOSEQUENZA DI UNA SEQUENZA DATA  $X$  E' UNA SEQUENZA OTTENUTA CANCELLANDO DA  $X$  ZERO O PIU' ELEMENTI (MANTENENDO L'ORDINE).

PIU' PRECISAMENTE, SE  $X = x_1 x_2 \dots x_m$ , UNA SEQUENZA  $Z = z_1 z_2 \dots z_k$  E' UNA SOTTOSEQUENZA DI  $X$  SE ESISTE UNA SEQUENZA  $\langle i_1, i_2, \dots, i_k \rangle$  DI INDICI TALE CHE:

$$- 1 \leq i_1 < i_2 < \dots < i_k \leq m$$

$$- z_j = x_{i_j}, \text{ PER OGNI } j = 1, 2, \dots, k$$

$$X = \dots x_{i_1} \dots x_{i_2} \dots x_{i_{k-1}} \dots x_{i_k} \dots$$
$$Z = z_1 \quad z_2 \quad \dots \quad z_{k-1} \quad z_k$$

ESEMPIO SIA  $X = \overset{1}{A} \overset{2}{B} \overset{3}{C} \overset{4}{B} \overset{5}{D} \overset{6}{A} \overset{7}{B}$ .

ALLORA  $Z = \underset{1}{B} \underset{2}{C} \underset{3}{D} \underset{4}{B}$  E' UNA SOTTOSEQUENZA DI  $X$

CORRISPONDENTE ALLA SEQUENZA DI INDICI  $\langle 2, 3, 5, 7 \rangle$

DEF DATE DUE SEQUENZE  $X$  E  $Y$ , DICIAMO CHE  $Z$  E' UNA SOTTOSEQUENZA COMUNE DI  $X$  E  $Y$  SE  $Z$  E' UNA SOTTOSEQUENZA DI ENTRAMBE LE SEQUENZE  $X$  E  $Y$

ESEMPIO DATE  $X = \overset{1}{A} \overset{2}{B} \overset{3}{C} \overset{4}{B} \overset{5}{D} \overset{6}{A} \overset{7}{B}$  E  
 $Y = B D C A B A$  ,

- LA SEQUENZA  $Z = B C A$  E' UNA SOTTOSEQUENZA COMUNE DI  $X$  E  $Y$ .
- TUTTAVIA  $B C A$  NON E' LA PIU' LUNGA SOTTOSEQUENZA COMUNE DI  $X$  E  $Y$ , IN QUANTO  $B C B A$  ,  $B C A B$  ,  $B D A B$  SONO SOTTOSEQUENZE COMUNI DI  $X$  E  $Y$  (DI LUNGHEZZA MASSIMA).

## PROBLEMA DELLA PIÙ LUNGA SOTTOSEQUENZA COMUNE:

DATE DUE SEQUENZE  $X$  E  $Y$  DETERMINARE UNA SOTTOSEQUENZA DI LUNGHEZZA MASSIMA ( $LCS$ ) CHE SIA COMUNE A  $X$  E  $Y$ .

1 2 3 4 5 6

## SOLUZIONE MEDIANTE RICERCA ESAUSTIVA

È ESPONENZIALE IN  $\min(|X|, |Y|)$ , IN QUANTO UNA SEQUENZA DI LUNGHEZZA  $m$  HA ESATTAMENTE  $2^m$  SOTTOSEQUENZE.

COMPLESSITÀ:  $O(\max(|X|, |Y|) \cdot 2^{\min(|X|, |Y|)})$

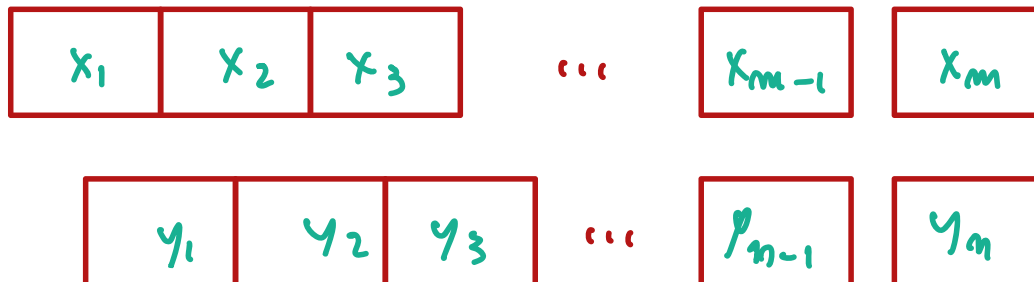
- IL PROBLEMA DELLA  $LCS$  PUÒ ESSERE RISOLTO IN MODO EFFICIENTE UTILIZZANDO LA PROGRAMMAZIONE DINAMICA.

# FASE 1: CARATTERIZZAZIONE DELLA PIÙ LUNGA SOTTOSEQUENZA COMUNE

NOTAZIONE DATA UNA SEQUENZA  $X = x_1 x_2 \dots x_m$ ,  
PONIAMO  $X_i := x_1 x_2 \dots x_i$ , PER  $i = 0, 1, 2, \dots, m$ .

INOLTRE PER OGNI SIMBOLO  $a$  DELL'ALFABETO  
PONIAMO  $Xa := x_1 x_2 \dots x_m a$

ESEMPIO: SE  $X = A B C B D A B$ , ALLORA  
 $X_1 = A$ ,  $X_4 = A B C B$ ,  $X_0 = \varepsilon$  (SEQUENZA VUOTA)  
 $X_3 D = A B C D$ , ECC.



# TEOREMA (SOTTOSTRUTTURA OTTIMA DI UNA LCS)

SIANO DATE DUE SEQUENZE  $X = x_1 x_2 \dots x_m$  E

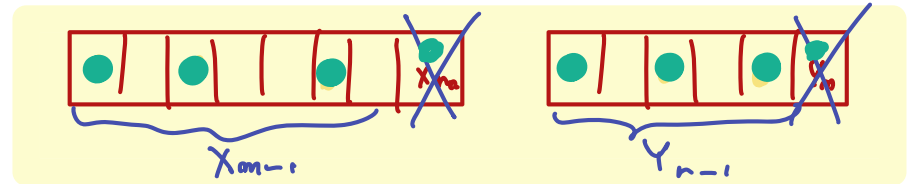
$Y = y_1 y_2 \dots y_n$  TALI CHE  $m \geq 1$  E  $n \geq 1$ ,

E SIA  $Z = z_1 z_2 \dots z_k$  UNA LCS DI  $X$  E  $Y$ .

1. SE  $x_m = y_m$ , ALLORA

-  $z_k = x_m = y_m$  E

-  $Z_{k-1}$  E' UNA LCS DI  $X_{m-1}$  E  $Y_{n-1}$

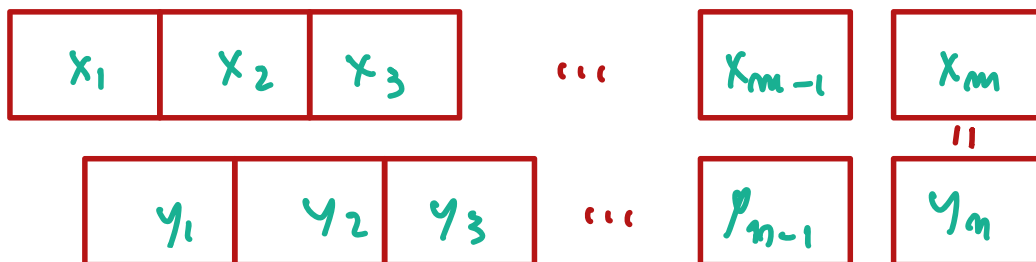


2. SE  $x_m \neq y_m$ , ALLORA

$z_k \neq x_m \implies Z$  E' UNA LCS DI  $X_{m-1}$  E  $Y$

3. SE  $x_m \neq y_m$ , ALLORA

$z_k \neq y_m \implies Z$  E' UNA LCS DI  $X$  E  $Y_{n-1}$ .



DIM. (1) SE  $z_k \neq x_m$  ALLORA  $Z$  E' UNA  
SOTTOSEQUENZA COMUNE DI  $X_{m-1}$  E  $Y_{m-1}$   
E QUINDI  $Z_{x_m}$  E' UNA SOTTOSEQUENZA  
COMUNE DI  $X$  E  $Y$ , ASSURDO.

QUINDI  $z_k = x_m = y_m$  E PERTANTO  $Z_{k-1}$  E'  
UNA SOTTOSEQUENZA COMUNE DI  $X_{m-1}$  E  $Y_{m-1}$ .  
SE  $Z_{k-1}$  NON FOSSE DI LUNGHEZZA MASSIMA,  
ESISTEREBBE W LCS DI  $X_{m-1}$  E  $Y_{m-1}$  TALE  
CHE  $|W| > |Z_{k-1}|$ .

MA, ALLORA  $W_{x_m}$  SAREBBE UNA SOTTOSEQUENZA  
COMUNE DI  $X$  E  $Y$ , ASSURDO IN QUANTO  
 $|W_{x_m}| > |Z|$  E  $Z$  E' UNA LCS DI  
 $X$  E  $Y$

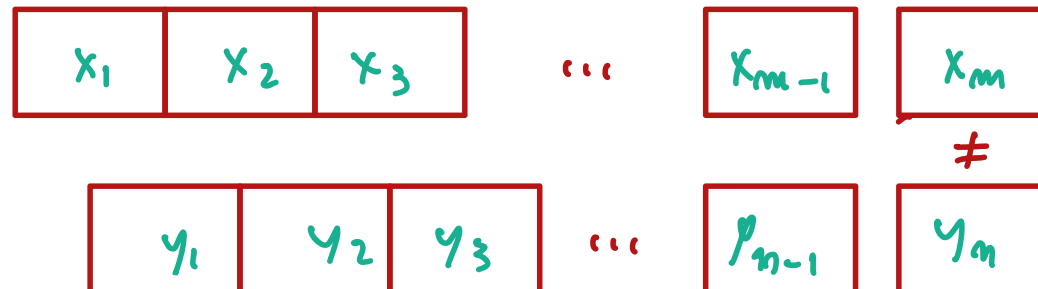


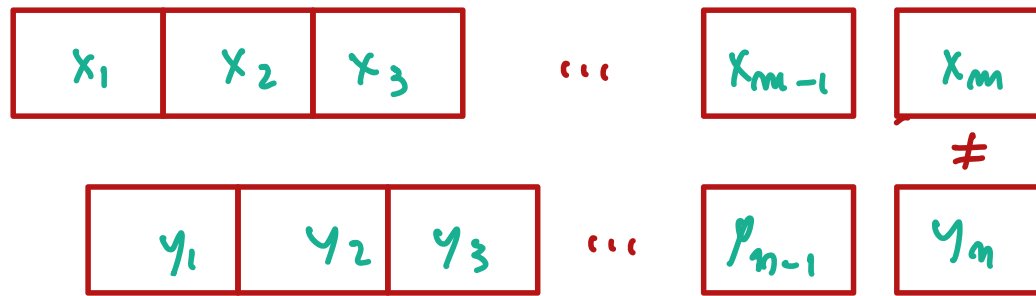
(2) SE  $z_k \neq x_m$  ALLORA  $Z$  E' UNA SOTTOSEQUENZA

$$x_m \neq y_m$$

COMUNE DI  $X_{m-1}$  E  $Y$ .

SE ESISTESSE UNA SOTTOSEQUENZA COMUNE  $W$  DI  $X_{m-1}$  E  $Y$  TALE CHE  $|W| > |Z|$ , CIO' SAREBBE ASSURDO IN QUANTO  $W$  SAREBBE A MAGGIOR RAGIONE UNA SOTTOSEQUENZA COMUNE DI  $X$  E  $Y$  (DI LUNGHEZZA MAGGIORE DI QUELLA DELLA LCS  $Z$ ).





(3) SE  $z_k \neq y_m$  ALLORA  $Z$  E' UNA SOTTOSEQUENZA  
 COMUNE DI  $X$  E  $Y_{m-1}$ .

$$x_m \neq y_m$$

SE ESISTESSE UNA SOTTOSEQUENZA COMUNE  $W$   
 DI  $X$  E  $Y_{m-1}$  TALE CHE  $|W| > |Z|$ ,  
 POICHE'  $W$  SAREBBE ANCHE UNA SOTTOSEQUENZA  
 COMUNE  $X$  E  $Y$ , VERREBBE CONTRADDETTA  
 LA MASSIMALITA' DI  $|Z|$ .

## SPAZIO DEI SOTTOPROBLEMI

$$(X, Y) = (X_m, Y_m) \mapsto (X_{m-1}, Y_{m-1}) / (X_m, Y_{m-1}) / (X_{m-1}, Y_m)$$

$$(X_{m-1}, Y_{m-1}) \mapsto (X_{m-2}, Y_{m-2}) / (X_{m-1}, Y_{m-2}) / (X_{m-2}, Y_{m-1})$$

$$(X_m, Y_{m-1}) \mapsto (X_{m-1}, Y_{m-2}) / (X_m, Y_{m-2}) / (X_{m-1}, Y_{m-1})$$

$$(X_{m-1}, Y_m) \mapsto (X_{m-2}, Y_{m-1}) / (X_{m-1}, Y_{m-1}) / (X_{m-2}, Y_m)$$

$\vdots$

$$\{(X_i, Y_j) : 0 \leq i \leq m, 0 \leq j \leq m\}$$

## SPAZIO DEI SOTTOPROBLEMI

DAL PRECEDENTE TEOREMA SEGUE CHE LO SPAZIO DEI SOTTOPROBLEMI E'  $\{(X_i, Y_j) : 0 \leq i \leq m, 0 \leq j \leq m\}$  LA CUI CARDINALITA' E'  $O(m^2)$ .

## FASE 2: SOLUZIONE RICORSIVA

DEFINIAMO  $c[i, j]$ , PER  $0 \leq i \leq m$  E  $0 \leq j \leq m$ , COME LA LUNGHEZZA DI UNA LCS DI  $X_i$  E  $Y_j$ .  
IN VIRTU' DELLA SOTTOSTRUTTURA OTTIMA SI HA

$$c[i, j] = \begin{cases} 0 & \text{SE } i=0 \text{ O } j=0 \\ c[i-1, j-1] + 1 & \text{SE } i, j > 0 \text{ E } x_i = y_j \\ \max(c[i, j-1], c[i-1, j]) & \text{SE } i, j > 0 \text{ E } x_i \neq y_j \end{cases}$$

- NUMERO  $n_1$  DI SOTTOPROBLEMI UTILIZZATI IN UNA SOLUZIONE OTTIMA  $= 1$
- NUMERO  $n_2$  DI SCELTE PER DETERMINARE QUALI SOTTOPROBLEMI UTILIZZARE  $\leq 2$

### FASE 3: CALCOLO DELLA LUNGHEZZA DI UNA LCS

LCS\_LENGTH( $X, Y$ )

$m := \text{length}[X]$

$n := \text{length}[Y]$

for  $i := 1$  to  $m$  do

$c[i, 0] := 0$

for  $j := 0$  to  $n$  do

$c[0, j] := 0$

CASO BASE

$O(m+n)$

$c[i, j]$

	0	1	2	3	4	5
0	0	0	0	0	0	0
1	0	>	>	>	>	>
2	0	>	>	>	>	>
3	0	>	>	>	>	>
4	0	>	>	>	>	>

for  $i := 1$  to  $m$  do

for  $j := 1$  to  $n$  do

if  $x_i = y_j$  then

$c[i, j] := c[i-1, j-1] + 1$  ;  $b[i, j] := "\nwarrow"$

else if  $c[i-1, j] \geq c[i, j-1]$  then

$c[i, j] := c[i-1, j]$  ;  $b[i, j] := "\uparrow"$

else  $c[i, j] := c[i, j-1]$  ;  $b[i, j] := "\leftarrow"$

return  $c, b$

DEFINIZIONE  
RICORSIVA

$O(mn)$

# COMPLESSITA' DI LCS\_LENGTH: $\Theta(mn)$

ESEMPIO: DATI  $X = A B C B D A B$   
 $Y = B D C A B A$ .

LCS\_LENGTH( $X, Y$ ) CALCOLA LA TABELLA.

		0	1	2	3	4	5	6
	$y_j$	B	D	C	A	B	A	
0	$x_i$	0	0	0	0	0	0	0
1	A	0	0	0	0	1	1	1
2	B	0	1	1	1	2	2	2
3	C	0	1	1	2	2	2	2
4	B	0	1	1	2	3	3	3
5	D	0	1	2	2	2	3	3
6	A	0	1	2	2	3	3	4
7	B	0	1	2	2	3	4	4

DA QUESTA SI  
EVINCE SUBITO CHE  
LA LUNGHEZZA DI UNA  
LCS E' 4

$Z =$  B C B A  
B C A B  
B D A B

## FASE 4: COSTRUZIONE DI UNA LCS

PRINT-LCS( $b, X, i, j$ )

if  $i=0$  o  $j=0$  then  
return

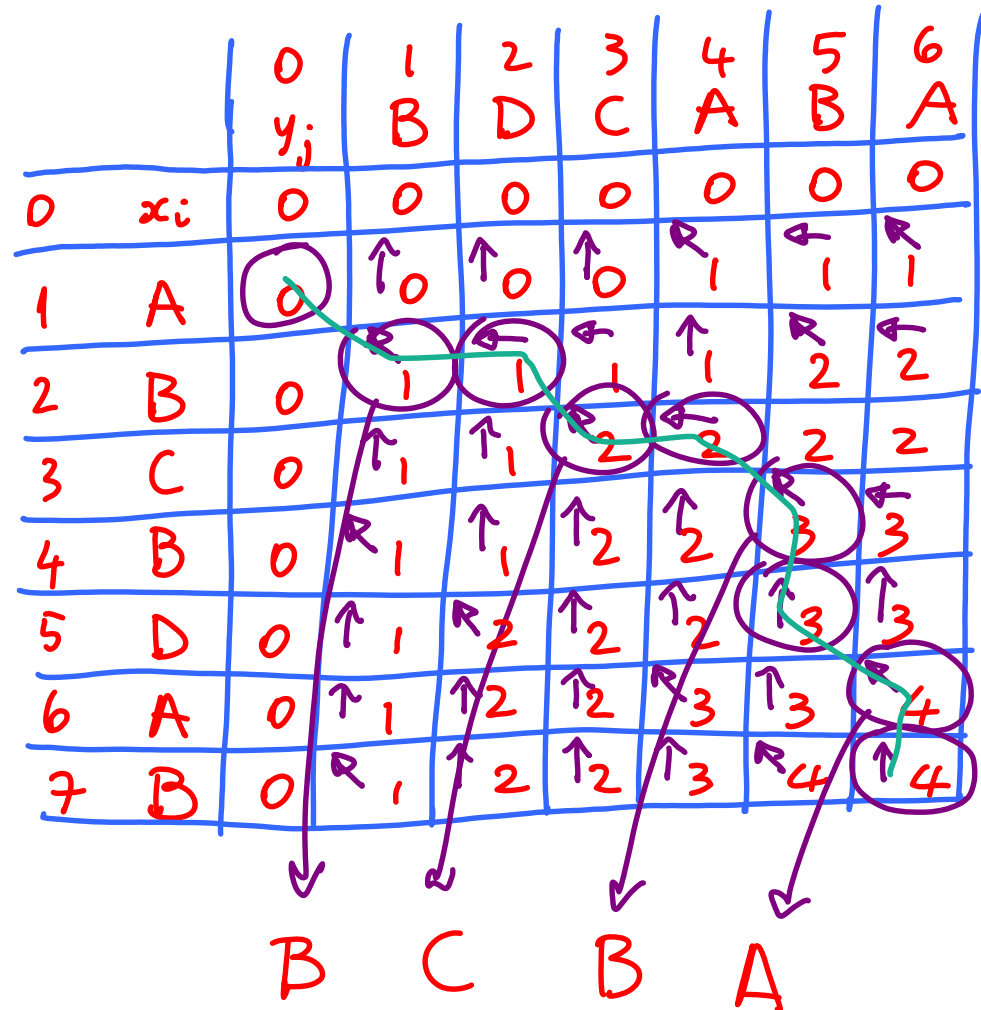
if  $b[i, j] = "\nwarrow"$  then  
PRINT-LCS( $b, X, i-1, j-1$ )  
stampa  $x_i$

elseif  $b[i, j] = "\uparrow"$  then  
PRINT-LCS( $b, X, i-1, j$ )

else  
PRINT-LCS( $b, X, i, j-1$ )

COMPLESSITA':  $O(m+n)$

PRINT-LCS( $b, X, m, m$ )



E' UNA LCS

## ESERCIZI

15.4-4 SPIEGARE COME CALCOLARE LA LUNGHEZZA DI UNA LCS UTILIZZANDO SOLTANTO 2  $\min(m, m)$  POSIZIONI NELLA TABELLA  $c$  PIÙ UNO SPAZIO  $O(1)$  AGGIUNTIVO, RISOLVERE LO STESSO PROBLEMA UTILIZZANDO  $\min(m, m)$  POSIZIONI PIÙ UNO SPAZIO  $O(1)$  AGGIUNTIVO,

15.4-5 PROGETTARE UN ALGORITMO  $O(n^2)$  PER TROVARE UNA PIÙ LUNGA SOTTOSEQUENZA CRESCENTE DI UNA SEQUENZA DI  $n$  NUMERI



IS.4-4

SPIEGARE COME CALCOLARE LA LUNGHEZZA DI UNA LCS  
 UTILIZZANDO SOLTANTO 2  $\min(m, n)$  POSIZIONI NELLA  
 TABELLA  $c$  PIÙ UNO SPAZIO  $O(1)$  AGGIUNTIVO,  
 RISOLVERE LO STESSO PROBLEMA UTILIZZANDO  $\min(m, n)$   
 POSIZIONI PIÙ UNO SPAZIO  $O(1)$  AGGIUNTIVO.

		B	D	C	A	B	A
A							
B							
C							
B							
D							
A							
B							

	0	0	0	0	0	0	0
0	0	0	0	1	1	1	
0	1	1	1	1	2	2	
0	1	1	2	2	2	2	
0	1	1	2	2	3	3	
0	1	2	2	2	3	3	
0	1	2	2	3	3	4	
0	1	2	2	3	4	4	

IS.4-4

SPIEGARE COME CALCOLARE LA LUNGHEZZA DI UNA LCS  
 UTILIZZANDO SOLTANTO 2  $\min(m, n)$  POSIZIONI NELLA  
 TABELLA C PIÙ UNO SPAZIO  $O(1)$  AGGIUNTIVO,  
 RISOLVERE LO STESSO PROBLEMA UTILIZZANDO  $\min(m, n)$   
 POSIZIONI PIÙ UNO SPAZIO  $O(1)$  AGGIUNTIVO,

		B	D	C	A	B	A
A	0	1	1	1	1	2	2
B							
C							
B							
D							
A							
B							
	0	1	1	2	2	2	2

0	0	0	0	1	1	1
0	1	1	1	1	2	2
0	1	1	2	2	2	2
0	1	1	1	1	3	3
0	1	2	2	2	3	3

0	1	2	2	3	4	4
---	---	---	---	---	---	---

15.4-5 PROGETTARE UN ALGORITMO  $O(n^2)$  PER TROVARE UNA PIÙ LUNGA SOTTOSEQUENZA CRESCENTE DI UNA SEQUENZA DI  $n$  NUMERI

SIA  $S$  UNA SEQUENZA DI NUMERI

- ORDINARE LA SEQUENZA  $S$  ELIMINANDO LE RIPETIZIONI.

SIA  $S_1$  LA SEQUENZA OTTENUTA

- CERCARE LA LCS DI  $S$  ED  $S_1$ .

COMPLESSITA':  $O(n^2)$