

# Longest common subsequence

Given two sequences  $x = [1, \dots, m]$  and  $y = [1, \dots, n]$  determine ② longest common subsequence

$LCS(x, y)$   $\rightsquigarrow$  functional notation even though it's not a function

Ex.

$x: A B C B D A B \rightarrow \begin{cases} B D A B \\ B C A B \\ B C B A \end{cases}$   
 $y: B D C A B A$

Brute force algorithm

① Check every subsequence of  $x[1 \dots m]$  to see if it is also a sequence of  $y[1 \dots n]$

Analysis: 1. to check if a sequence is a subsequence of  $y$   
 $\hookrightarrow O(n) \rightarrow$  you have to get the first char matching and then you scan

2. How many subsequences in  $x$ ?

$\hookrightarrow 2^m$

$\Rightarrow O(n \cdot 2^m)$

## Simplification Step

1. Look at the length of  $LCS(x, y) \rightsquigarrow c[x, y]$
2. Extend 1. to return the actual LCS

Strategy: Consider only prefixes of  $x$  and  $y$  and we are going to show how we can express the length of the LCS of prefixes in terms of each other.

Define:  $c[i, j] = |LCS(x[1 \dots i], y[1 \dots j])|$

$\hookrightarrow$  we need to calculate  $c[i, j]$  (or  $c_{ij}$ )  $\forall i, j \in n, m$

How do we solve the problem of LCS(x, y) if we have  $c_{ij} \forall i, j$ ?

$$L \rightarrow c[m, n] = |LCS(x, y)|$$

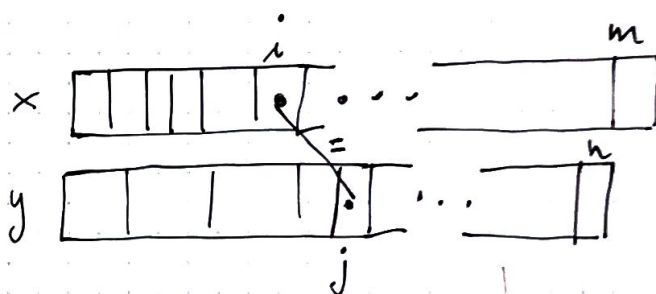
We want to express the general  $c[m, n]$  in terms of other  $c[i, j]$

Theorem

$$c[i, j] = \begin{cases} 0, & \text{if } i=0 \text{ or } j=0 \\ c[i-1, j-1] + 1, & \text{if } x[i] = y[j] \\ \max \{ c[i, j-1], c[i-1, j] \}, & \text{otherwise} \end{cases}$$

Proof.

case  $x[i] = y[j]$



Let  $z[1 \dots k] = LCS(x[1 \dots i], y[1 \dots j])$  where  $c[i, j] = k$

Then,  $z[k] = x[i] = y[j]$  because if  $z$  does not include  $x[i]$  or  $y[j]$  then we can add this char to  $z$  and make it longer because  $x[i] = y[j] \rightarrow$  it can add an element to  $z' = z + x[i]$  then  $z$  could not be a LCS...

Thus  $z[1 \dots k-1]$  is a CS of  $x[1 \dots i-1]$  and  $y[1 \dots j-1]$

CLAIM:  $z[1 \dots k-1]$  is a LCS of  $x[1 \dots i-1]$  and  $y[1 \dots j-1]$

Proof of the claim: Absurd to suppose that  $w$  is a longer CS than  $z \Rightarrow |w| > k-1$

We use the CUT & PASTE argument

$w \parallel z[k] \xrightarrow{\text{string concatenation}}$  is certainly a CS of  $x[1 \dots i]$  and  $y[1 \dots j]$  and has length greater than  $k$   
 $\rightarrow$  CONTRADICTION

Thus,  $c[i-1, j-1] = k-1$ , which implies that

$$c[i, j] = c[i-1, j-1] + 1.$$

The other case can be proven in the same way.  $\square$

We have an OPTIMAL SUBSTRUCTURE because the optimal solution of the subproblems compose the solution of the whole problem.

↳ If  $z = \text{LCS}(x, y)$  then any prefix of  $z$  is a LCS of a prefix of  $x$  and a prefix of  $y$ .

Now we have a strategy to compute LCS  
RECURSIVE algorithm

$\text{LCS}(x, y, i, j)$

// we ignore the base cases  $i=0, j=0$

if  $x[i] == y[j]$  then

$$c[i, j] = \text{LCS}(x, y, i-1, j-1) + 1$$

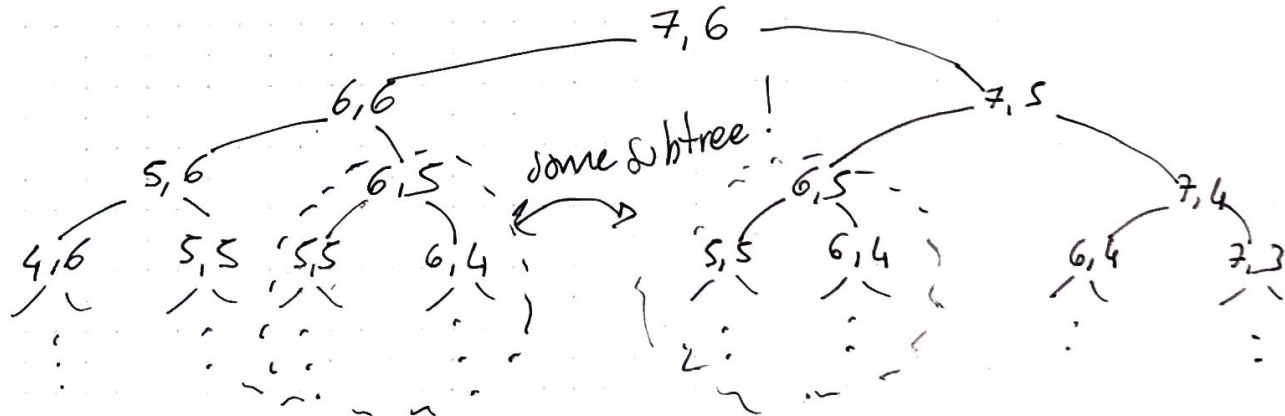
else

$$c[i, j] = \max(\text{LCS}(x, y, i-1, j), \text{LCS}(x, y, i, j-1))$$

return  $c[i, j]$

worst-case analysis: the second clause is problematic because there are 2 LCS calls

RECURSION TREE for  $n=7, m=6$





The height of the tree is  $m+n$  so being a binary tree each level brings  $2^i$  work  $\approx O(2^{m+n})$

$\hookrightarrow$  i.e. LCS contains  $m \cdot n$  distinct subproblems.

$$T(n, m) = \begin{cases} T(n-1, m-1) + O(1) & \text{if } x[n] = y[m] \\ T(n-1, m) + T(n, m-1) + O(1), & \text{otherwise} \end{cases}$$

$$T(n, m) = 2^{n-1} T(0, m) + \dots + \dots \quad \text{if } m < n \quad (\text{longest branch of height } n)$$

$\downarrow$   
worst-case

$$T(n, m) = 2^{m-1} T(0, n) + \dots \quad \text{if } n < m$$

BOTTOM-UP. MEMOIZED ALGORITHM

LCS-LENGTH( $x, y, m, n$ )

let  $b[1..m, 1..n]$  and  $c[0..m, 0..n]$  be new tables

for  $i = 1$  to  $m$  do

$c[i, 0] = 0$

for  $j = 0$  to  $n$  do

$c[0, j] = 0$

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

$c[i, j] = c[i, j-1]$

$b[i, j] = \rightarrow$

return  $c, b$

Example

n

j	0	1	2	3	4	5	6
i \ y <sub>j</sub>	B	D	C	A	B	A	
0 x <sub>i</sub>	0	0	0	0	0	0	0
1 A	0	↑	↑	↑	↖	←	↖
2 B	0	↖	←	←	↑	↖	←
3 C	0	↑	↑	↖	←	↑	↑
4 D	0	↖	↑	↑	↑	↖	←
5 D	0	↑	↖	↑	↑	↑	↖
6 A	0	↑	↑	↑	↖	↑	↖
m=7 B	0	↖	↑	↑	↑	↖	↑

⇒ D C B A

PRINT-LCS(b, x, i, j)

if  $i == 0$  or  $j == 0$   
return m

if  $b[i, j] == "↖"$  then

PRINT-LCS(b, x, i-1, j-1)

print x<sub>i</sub>

else if  $b[i, j] == "↑"$  then

PRINT-LCS(b, x, i-1, j) // we go up

else PRINT-LCS(b, x, i, j-1) // we go left