# Dynamic programming

Longest Common Subsequence

## Dynamic programming

- It is used, when the solution can be recursively described in terms of solutions to subproblems (optimal substructure)
- Algorithm finds solutions to subproblems and stores them in memory for later use
- More efficient than "brute-force methods", which solve the same subproblems over and over again

#### Longest Common Subsequence (LCS)

**Subsequence:** Subset of letters in order from left to right.

Example:  $X = \{A B C B D A B \}, Y = \{B D C A B A\}.$ 

#### **Remember:**

- "Subsequence" is different from "substring". Substring is continuous, no gap. But, subsequence can have gap.
- Example, CBD is a substring of X. BBA is a subsequence of X, but not a substring of X. DBBB is not a subsequence of X, because order violation.
- Example: BDAB is both substring and subsequence of X, But it is only a subsequence of Y
- Substrings are also subsequence, but a subsequence may not be a substring.

**Longest Common Subsequence:** Same subsequence of maximum length.

- Example: X = A B C B D A B, Y = B D C A B A = BCBA (4)
- Other examples: BCAB(4), BDAB(4). Maximum possible length = 4

#### Why LCS Algorithm (important)

#### Brute Force Algorithm:

- if |X| = m, |Y| = n, then there are  $2^m$  subsequences of x; we must compare each with each subsequence of Y  $(2^n)$
- So the running time of brute-force algorithm is  $\geq O(2^n 2^m)$
- For example, X=AMP, Y=AP.
  - All possible subsequence of X are:  $\phi$ , A, M, P, AM, AP, MP AMP = total  $2^3=8$
  - For Y, they are:  $\phi$ , A, P, AP = total  $2^2=4$
  - Compare each of X with each of Y = total compare  $2^{3*}2^{2} = 32$
- But note that: there are many repetitions. For example, A of X and A of Y are compared 4 times (in A, AM, AP, AMP of X)\*2 times (in A, AP of Y) = 8 times.
- So we can save time, by saving this comparison into a table. This is called *optimal substructure property of dynamic programming*
- Actually, we will save LCS of prefixes of X and Y

#### LCS Algorithm

- First we'll find the length of LCS. Later we'll modify the algorithm to find LCS itself.
- Define  $X_i$ ,  $Y_j$  to be the prefixes of X and Y of length i and j respectively
- Define c[i,j] to be the length of LCS of  $X_i$  and  $Y_j$
- Then the length of LCS of X and Y will be c[m,n]

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

#### LCS recursive solution

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

- We start with i = j = 0 (empty substrings of x and y)
- Since  $X_0$  and  $Y_0$  are empty strings, their LCS is always empty (i.e. c[0,0] = 0)
- LCS of empty string and any other string is empty, so for every i and j: c[0, j] = c[i, 0] = 0

#### LCS recursive solution

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

- When we calculate c[i,j], we consider two cases:
- First case: x[i]=y[j]: one more symbol in strings X and Y matches, so the length of LCS  $X_i$  and  $Y_j$  equals to the length of LCS of smaller strings  $X_{i-1}$  and  $Y_{i-1}$ , plus 1

#### LCS recursive solution

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

- Second case: x[i] != y[j]
- As symbols don't match, our solution is not improved, and the length of LCS(X<sub>i</sub>, Y<sub>j</sub>) is the same as before (i.e. maximum of LCS(X<sub>i</sub>, Y<sub>j-1</sub>) and LCS(X<sub>i-1</sub>,Y<sub>j</sub>)

Why not just take the length of LCS( $X_{i-1}$ ,  $Y_{j-1}$ ) ?

#### LCS Length Algorithm

```
LCS-Length(X, Y)
1. m = length(X) // get the # of symbols in X
2. n = length(Y) // get the # of symbols in Y
3. for i = 1 to m c[i,0] = 0 // special case: Y_0
4. for j = 1 to n c[0,j] = 0 // special case: X_0
5. for i = 1 to m
                               // for all X<sub>i</sub>
6. for j = 1 to n
                                      // for all Y<sub>i</sub>
7.
             if (X_i == Y_i)
8.
                   c[i,j] = c[i-1,j-1] + 1
             else c[i,j] = max(c[i-1,j], c[i,j-1])
```

10. return c

9

#### LCS Example

We'll see how LCS algorithm works on the following example:

- $\mathbf{X} = \mathbf{ABCB}$
- $\mathbf{Y} = \mathbf{BDCAB}$

What is the Longest Common Subsequence of X and Y?

$$LCS(X, Y) = BCB$$
  
 $X = A B C B$   
 $Y = B D C A B$ 

## LCS Example (0)

ABCB BDCAB

	j	0	1	2	3	4	5 <sup>L</sup>
i		Yj	В	D	C	A	В
0	Xi						
1	A						
2	В						
3	C						
4	В						

$$X = ABCB$$
;  $m = |X| = 4$   
 $Y = BDCAB$ ;  $n = |Y| = 5$   
Allocate array c[5,4]

## LCS Example (1)

ABCB BDCAB

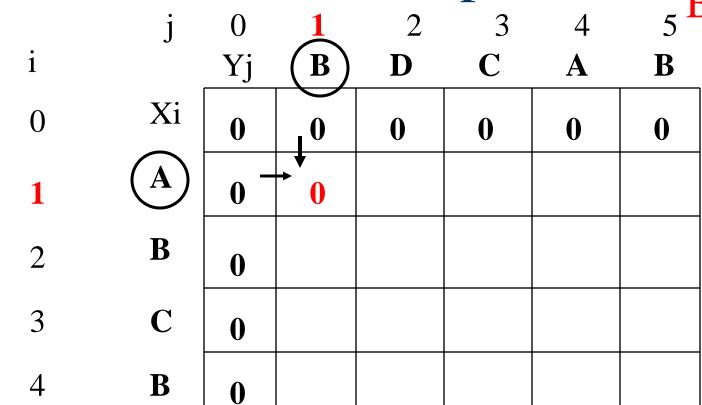
i		Yj	B	D	$\mathbf{C}$	$\mathbf{A}$	B
	J	U	1	2	3	4	3

		<i>J</i>					
0	Xi	0	0	0	0	0	0
1	A	0					
2	В	0					
3	C	0					
4	В	0					

for 
$$i = 1$$
 to m  $c[i,0] = 0$   
for  $j = 1$  to n  $c[0,j] = 0$ 

## LCS Example (2)

RDCAR



if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

## LCS Example (3)

RDCAR

	j	0	1	2	3	4	5
i		Yj	В	D	C	A	В
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0		
2	В	0					
3	C						

 $\mathbf{B}$ 

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

## LCS Example (4)

RDC A R

					_		
	j	0	1	2	3	4	5 E
i		Yj	B	$\mathbf{D}$	$\mathbf{C}$	$(\mathbf{A})$	В
0	Xi	0	0	0	0 、	0	0
1	(A)	0	0	0	0	1	
2	В	0					
3	C	0					
4	В	0					

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

## LCS Example (5)

RDC A R

							F
	j	0	1	2	3	4	5
i		Yj	В	D	C	A	(B)
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1 -	<b>1</b>
2	В	0					
3	C	0					
4	В	0					

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

## LCS Example (6)

ABCB BDCAR

						•	<b>—</b>
	j	0	1	2	3	4	5
i		Yj	<b>(B)</b>	D	$\mathbf{C}$	A	В
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	B	0	1				
3	C	0					
4	В	0					

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

#### LCS Example (7)

ABCB BDCAB

	j	0	1	2	3	4	5
i	_	Yj	В	D	C	A	<b>B</b>
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	$egin{pmatrix} \mathbf{B} \end{pmatrix}$	0	1	1	1	<b>→</b> 1	
3	C	0					
4	В	0					

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

# LCS Example (8)

ABCB RDCAR

							$\mathbf{p}$
	j	0	1	2	3	4	<b>5</b> D
i		Yj	$\mathbf{B}$	D	$\mathbf{C}$	A	(B)
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1,	1
2	B	0	1	1	1	1	2
3	C	0					

 $\mathbf{B}$ 

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

### LCS Example (10)

RDCAR

	j	0	1	_2	3	4	5
i		Yj	B	D)	C	A	В
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	В	0	1	_1	1	1	2
3	$\bigcirc$	0	<b>1</b> -	<b>1</b>			
4	В	0					

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

## LCS Example (11)

RDCAR

	j	0	1	2	3	4	5
i		Yj	В	D	(C)	A	В
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	В	0	1	1,	1	1	2
3	$\bigcirc$	0	1	1	2		
4	В	0					

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

# LCS Example (12)

Yj B D B Xi B B

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

## LCS Example (13)

RDC A R

	j	0	1	2	3	4	5 <sup>D</sup>
i		Yj	B	D	C	A	В
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	В	0	1	1	1	1	2
3	C	0 🔪	1	1	2	2	2
4	B	0	1				

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

## LCS Example (14)

ABCB BDCAB

	j	0	1	2	3	4	5 <sup>E</sup>
i		Yj	В	D	C	A	<b>S</b> B
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	В	0	1	1	1	1	2
3	C	0	1	1	_2	2	2
4	$oxed{B}$	0	1 -	<b>1</b>	<sup>†</sup> <sub>2</sub> -	<b>2</b>	

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

#### LCS Example (15)

RDCAR

	j	0	1	2	3	4	5	
i		Yj	В	D	C	A	B	
0	Xi	0	0	0	0	0	0	
1	A	0	0	0	0	1	1	
2	В	0	1	1	1	1	2	
3	C	0	1	1	2	2 🔨	2	
4	lacksquare	O	1	1	2	2	3	

if 
$$(X_i == Y_j)$$
  
 $c[i,j] = c[i-1,j-1] + 1$   
else  $c[i,j] = max(c[i-1,j],c[i,j-1])$ 

### LCS Algorithm Running Time

- LCS algorithm calculates the values of each entry of the array c[m,n]
- So what is the running time?

O(m\*n)

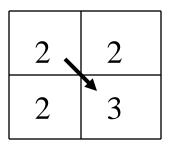
since each c[i,j] is calculated in constant time, and there are m\*n elements in the array

#### How to find actual LCS

- So far, we have just found the *length* of LCS, but not LCS itself.
- We want to modify this algorithm to make it output Longest Common Subsequence of X and Y

Each c[i,j] depends on c[i-1,j] and c[i,j-1] or c[i-1,j-1]

For each c[i,j] we can say how it was acquired:



For example, here c[i,j] = c[i-1,j-1] + 1 = 2+1=3

#### How to find actual LCS - continued

Remember that

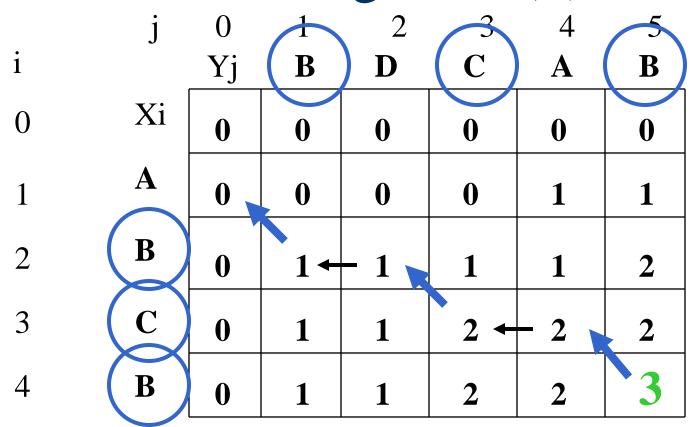
$$c[i, j] = \begin{cases} c[i-1, j-1] + 1 & \text{if } x[i] = y[j], \\ \max(c[i, j-1], c[i-1, j]) & \text{otherwise} \end{cases}$$

- So we can start from c[m,n] and go backwards by **following** arrow in reverse order (Step 1)
- Whenever c[i,j] = c[i-1,j-1]+1, remember x[i] (Step 2) (because x[i] is a part of LCS)
- When i=0 or j=0 (i.e. we reached the beginning), **output** remembered letters in reverse order (Final step)
- For each reverse path, there is one LCS. In this way all LCS of maximum length can be found.

# Finding LCS

	j	0	1	2	3	4	5
i		Yj	В	D	C	A	В
0	Xi	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	В	0	<b>1</b> ←	- 1 ×	1	1	2
3	C	0	1	1	2 +	- 2	2
4	В	0	1	1	2	2	3

## Finding LCS (2)



LCS (reversed order): B C B

LCS (straight order):

B C B

(this string turned out to be a palindrome)

#### LCS Example

We'll see how LCS algorithm works on the following example:

- $\mathbf{X} = \mathbf{PMDX}$
- $\mathbf{Y} = \mathbf{MPXD}$

What is the Longest Common Subsequence of X and Y?

$$LCS(X, Y) = MX, PD$$

See the video for this example.