# Dynamic Programming Longest Common Subsequence

CSCI 3100

---

## But first …

Follow up on the "sorting challenge" from last class

Which algorithm should we use to sort an "almost in order" array

Narrowed down to two options:
◦ Selection sort
◦ Merge sort

# Dynamic Programming

An algorithm design technique similar to divide and conquer but unlike divide&conquer, subproblems may overlap in this case.

Divide and conquer
- Partition the problem into subproblems (may overlap)
- Solve the subproblems recursively
- Combine the solutions to solve the original problem

Used for **optimization problems**

- Goal: **find an optimal solution** (minimum or maximum)

- There may be many solutions that lead to an optimal value

# Dynamic Programming

Applicable when subproblems are **not** independent

- Subproblems share subsubproblems

*e.g.:* Combinations:

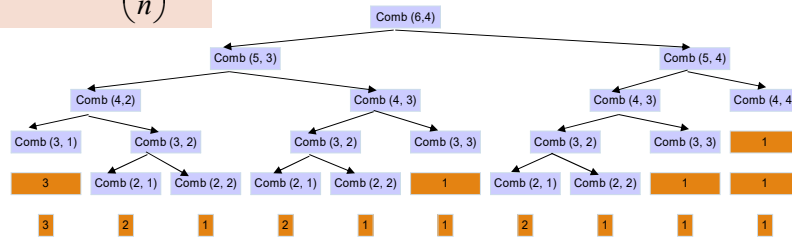$$\binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1}$$

$$\binom{n}{1} = n \qquad \binom{n}{n} = 1$$

- Dynamic programming solves every subproblem and stores the answer in a table

## Example: Combinations

$$\binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1}$$

$$\binom{n}{1} = n \qquad \binom{n}{n} = 1$$



# Dynamic Programming Algorithm

**Characterize** the structure of an optimal solution

**Recursively** define the value of an optimal solution
◦ An optimal solution to a problem contains within it an optimal solution to subproblems.
◦ Typically, the recursion tree contains many overlapping subproblems

**Compute** the value of an optimal solution in a bottom-up fashion
◦ Optimal solution to the entire problem is build in a bottom-up manner from optimal solutions to subproblems

**Construct** an optimal solution from computed information

# Longest Common Subsequence

Given two sequences

$$X = \langle x_1, x_2, \ldots, x_m \rangle$$

$$Y = \langle y_1, y_2, \ldots, y_n \rangle$$

find a **maximum length common subsequence** (LCS) of X and Y

*e.g.:*   If   $X = \langle A, B, C, B, D, A, B \rangle$

Subsequences of  X:

A subset of elements in the sequence taken in order

$\langle A, B, D \rangle, \langle B, C, D, B \rangle, \langle B, C, D, A, B \rangle$  etc.

# Example

$X = \langle A, B, C, B, D, A, B \rangle$      $X = \langle A, B, C, B, D, A, B \rangle$

$Y = \langle B, D, C, A, B, A \rangle$          $Y = \langle B, D, C, A, B, A \rangle$

$\langle \mathbf{B, C, B, A} \rangle$ and $\langle \mathbf{B, D, A, B} \rangle$ are
**longest common subsequences** of X and Y   (*length* = 4)

$\langle B, C, A \rangle$, however, is not a LCS of X and Y

8

# Applications of LCS

Molecular biology
- ◦ DNA sequences represented as combinations of letters ACGT
- ◦ Find how similar two sequences are

File comparison:
- ◦ Linux 'diff' command to compare two files

# Brute-Force Solution

For every subsequence of X, check whether it's a subsequence of Y
- ◦ There are $2^m$ subsequences of X to check

Each subsequence takes $\Theta(n)$ time to check
- ◦ scan Y for first letter, from there scan for second, and so on

**Running time**:   $\Theta(n2^m)$

10

## Making the choice

$X = \langle A, B, D, G, E \rangle$

$Y = \langle Z, B, D, E \rangle$

**Choice**: include one element into the common sequence (E) and solve the resulting subproblem

$X = \langle A, B, D, G \rangle$

$Y = \langle Z, B, D \rangle$

**Choice**: exclude an element from a string and solve the resulting subproblem

11

## Notations

Given a sequence $X = \langle x_1, x_2, ..., x_m \rangle$

we define the **i-th prefix** of X, for i = 0, 1, 2, …, m

$$X_i = \langle x_1, x_2, ..., x_i \rangle$$

**c[i, j]** = the **length** of a LCS of the sequences

$X_i = \langle x_1, x_2, ..., x_i \rangle$ and $Y_j = \langle y_1, y_2, ..., y_j \rangle$

12

# A Recursive Solution

**Case 1**: $x_i = y_j$

*e.g.:*     $X_i = \langle A, B, D, G, E \rangle$

$Y_j = \langle Z, B, D, E \rangle$

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

- Append $x_i = y_j$ to the LCS of $X_{i-1}$ and $Y_{j-1}$
- Must find a LCS of $X_{i-1}$ and $Y_{j-1}$

# A Recursive Solution

**Case 2**:  $x_i \neq y_j$

*e.g.:*        $X_i = \langle A, B, D, G \rangle$

$Y_j = \langle Z, B, D \rangle$

- Must solve two problems
  - find a LCS of $X_{i-1}$ and $Y_j$:     $X_{i-1} = \langle A, B, D \rangle$ and $Y_j = \langle Z, B, D \rangle$
  - find a LCS of $X_i$ and $Y_{j-1}$:     $X_i = \langle A, B, D, G \rangle$ and $Y_{j-1} = \langle Z, B \rangle$

$$c[i, j] = max \{ c[i - 1, j], c[i, j-1] \}$$

Optimal solution to a problem includes optimal solutions to subproblems

14

# Overlapping Subproblems

To find a LCS of  $(X_m$ and $Y_n)$
- we may need to find the LCS between $X_m$ and $Y_{n-1}$ and that of $X_{m-1}$ and $Y_n$
- Both of the above subproblems has the subproblem of finding the LCS of $(X_{m-1}$ and $Y_{n-1})$

Subproblems share subsubproblems

15

# Computing the Length of the LCS

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



8

## Additional Information

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

**b & c:**

| | | 0 | 1 | 2 | 3 | | n |
|---|---|---|---|---|---|---|---|
| | $y_j$: | | A | C | D | | F |
| 0 | $x_i$ | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | A | 0 | | | | | |
| 2 | B | 0 | | | c[i-1,j] | | |
| 3 | C | 0 | c[i,j-1] | ↑ | | | |
| | | 0 | | | | | |
| m | D | 0 | | | | | |

j

A matrix $b[i, j]$:

- For a subproblem [i, j] it tells us what choice was made to obtain the optimal value

- If $x_i = y_j$
  $b[i, j] = " \searrow "$
- Else, if $c[i - 1, j] \geq c[i, j-1]$
  $b[i, j] = " \uparrow "$
  else
  $b[i, j] = " \leftarrow "$

---

## Example

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

$X = \langle A, B, C, B, D, A, B \rangle$

$Y = \langle B, D, C, A, B, A \rangle$

If $x_i = y_j$
  $b[i, j] = " \searrow "$
else if $c[i - 1, j] \geq c[i, j-1]$
  $b[i, j] = " \uparrow "$
else
  $b[i, j] = " \leftarrow "$

| | | 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 | ↑1 | ↖2 | ←2 |
| 3 | C | 0 | ↑1 | ↑1 | ↖2 | ←2 | ↑2 | ↑2 |
| 4 | B | 0 | ↖1 | ↑1 | ↑2 | ↑2 | ↖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 |

# Constructing a LCS

Start at $b[m, n]$ and follow the arrows

When we encounter a "↖" in $b[i, j] \Rightarrow x_i = y_j$ is an element of the LCS

|   |   | 0 | 1 B | 2 D | 3 C | 4 A | 5 B | 6 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 | ↑1 | ↖2 | ←2 |
| 3 | C | 0 | ↑1 | ↑1 | ↖(2) | ←(2) | ↑2 | ↑2 |
| 4 | B | 0 | ↖1 | ↑1 | ↑2 | ↑2 | ↖(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) |

# LCS-LENGTH(X, Y, m, n)

1.  **for** $i \leftarrow 1$ **to** $m$
2.     **do** $c[i, 0] \leftarrow 0$
3.  **for** $j \leftarrow 0$ **to** $n$
4.     **do** $c[0, j] \leftarrow 0$

If one of the sequences is empty, the length of the LCS is zero

5.  **for** $i \leftarrow 1$ **to** $m$
6.     **do for** $j \leftarrow 1$ **to** $n$
7.        **do if** $x_i = y_j$
8.           **then** $c[i, j] \leftarrow c[i - 1, j - 1] + 1$    Case 1:   $x_i = y_j$
9.               $b[i, j] \leftarrow$ "↖"
10.          **else if** $c[i - 1, j] \geq c[i, j - 1]$
11.             **then** $c[i, j] \leftarrow c[i - 1, j]$
12.                 $b[i, j] \leftarrow$ "↑"
13.             **else** $c[i, j] \leftarrow c[i, j - 1]$
14.                 $b[i, j] \leftarrow$ "←"
15. **return** $c$ and $b$

Case 2:   $x_i \neq y_j$

**Running time**:

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

1. **if** i = 0 or j = 0
2.  **then return**
3. **if** b[i, j] = "↖ "
4.    **then** PRINT-LCS(b, X, i – 1, j – 1)
5.       print $x_i$
6. **elseif** b[i, j] = "↑"
7.       **then** PRINT-LCS(b, X, i – 1, j)
8.       **else** PRINT-LCS(b, X, i, j – 1)

Initial call:   PRINT-LCS(b, X, length[X], length[Y])

**Running time**:

# Improving the Code

What can we say about how each entry $c[i, j]$ is computed?
- It depends only on $c[i -1, j - 1]$, $c[i - 1, j]$, and $c[i, j - 1]$
- Eliminate table b and compute in $O(1)$ which of the three values was used to compute $c[i, j]$
- We save $\Theta(mn)$ space from table b
- However, we do not asymptotically decrease the auxiliary space requirements: still need table c

If we only need the length of the LCS
- LCS-LENGTH works only on two rows of c at a time
  - The row being computed and the previous row
- We can reduce the asymptotic space requirements by storing only these two rows

22