Longest Common Subsequence (LCS)

C.S. Marron cmarron@umbc.edu

CMSC 441 — Algorithms

Outline

Sequences and Subsequences Definitions An Example

The LCS Problem Optimal Substructure

4 日 > 4 日 > 4 日 > 4 日 > 4 日 > 4 日 > 4 日 > 9 4 0 4

Sequences and Subsequences

Definitions

A sequence is an ordered multiset in which the elements are taken from an underlying set. The usual sequence notation uses subscripts, e.g.,

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

where x_i are elements of the underlying set.

1

A subsequence of a sequence is an ordered multi-subset. Using subscript notation, if $X = \langle x_1, x_2, \dots, x_n \rangle$ is a sequence, then a subsequence is

$$\langle x_{i_1}, x_{i_2}, \dots, x_{i_k} \rangle \quad \langle x_{3}, x_{7}, \dots, x_{21} \rangle$$
where $1 \leq i_1 < i_2 < \dots < i_k \leq n$.

Common Subsequences

Definitions

- If X and Y are sequences, Z is a common subsequence if it is a subsequence of both X and Y.
- A common subsequence of X and Y with maximum length is a longest common subsequences (LCS).

◆□▶ ◆□▶ ◆ ■▶ ◆ ■ ● のへで

An LCS is not necessarily unique, although its length is!

Set = EA, B, C, D]

- $\blacktriangleright X = \langle D, A, B, D, D, C, D, A \rangle$
- $\blacktriangleright Y = \langle C, B, A, B, D, C, A, D \rangle$
- \triangleright $\langle D, A, B \rangle$ is a subsequence of X, but *not* of Y.
- \$\langle D, A, D \rangle\$ is a subsequence of both X and Y. It is a common subsequence, but it is not a longest common subsequence.
- ► Z = (A, B, D, C, D) is a longest common subsequence of X and Y.

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● のへで

► Z' = (A, B, D, C, A) is also an LCS of X and Y. The sequence is not unique, but the length (5) is.

- $\blacktriangleright X = \langle D, A, B, D, D, C, D, A \rangle$
- $\blacktriangleright Y = \langle C, B, A, B, D, C, A, D \rangle$
- \triangleright $\langle D, A, B \rangle$ is a subsequence of X, but *not* of Y.
- \$\langle D, A, D \rangle\$ is a subsequence of both X and Y. It is a common subsequence, but it is not a longest common subsequence.
- Z = (A, B, D, C, D) is a longest common subsequence of X and Y.

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ の < @

Z' = ⟨A, B, D, C, A⟩ is also an LCS of X and Y. The sequence is not unique, but the length (5) is.

- $\blacktriangleright X = \langle \underline{D}, \underline{A}, B, D, \underline{D}, C, D, A \rangle$
- $\blacktriangleright Y = \langle C, B, A, B, D, C, A, D \rangle$
- \triangleright $\langle D, A, B \rangle$ is a subsequence of X, but *not* of Y.
- \$\langle D, A, D \rangle\$ is a subsequence of both X and Y. It is a common subsequence, but it is not a longest common subsequence.
- \triangleright Z = $\langle A, B, D, C, D \rangle$ is a *longest common subsequence* of X and Y.

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ の < @

► Z' = (A, B, D, C, A) is also an LCS of X and Y. The sequence is not unique, but the length (5) is.

- $\blacktriangleright X = \langle D, \underline{A}, \underline{B}, \underline{D}, D, \underline{C}, \underline{D}, A \rangle$
- $\blacktriangleright Y = \langle C, B, \underline{A}, \underline{B}, \underline{D}, \underline{C}, A, \underline{D} \rangle$
- \triangleright $\langle D, A, B \rangle$ is a subsequence of X, but *not* of Y.
- \$\langle D, A, D \rangle\$ is a subsequence of both X and Y. It is a common subsequence, but it is not a longest common subsequence.
- Z = (A, B, D, C, D) is a longest common subsequence of X and Y.

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ の < @

Z' = ⟨A, B, D, C, A⟩ is also an LCS of X and Y. The sequence is not unique, but the length (5) is.

- $\blacktriangleright X = \langle D, \overline{A}, \overline{B}, \overline{D}, D, \overline{C}, D, \overline{A} \rangle$
- $\blacktriangleright Y = \langle C, B, A, B, D, C, A, D \rangle$
- \triangleright $\langle D, A, B \rangle$ is a subsequence of X, but *not* of Y.
- \$\langle D, A, D \rangle\$ is a subsequence of both X and Y. It is a common subsequence, but it is not a longest common subsequence.
- $Z = \langle A, B, D, C, D \rangle \text{ is a longest common subsequence of } X$ and Y. Length 5

< □ ▶ < □ ▶ < □ ▶ < □ ▶ = □ ● ○ < ○

► Z' = (A, B, D, C, A) is also an LCS of X and Y. The sequence is not unique, but the length (5) is.

Optimization: find LCS length.

Using the Example

 $X = \langle D, A, B, D, D, C, D, A \rangle$. $Y = \langle C, B, A, B, D, C, A, D \rangle$. $Z = \langle A, B, D, C, D \rangle$ Look at what happens considering the last characters of each sequence.

What can we say when the last character of X and the last character of Y are different?

Optimal Substructure Theorem

Notation

If $X = \langle x_1, x_2, \dots, x_m \rangle$ is a sequence, then we define X_i , $1 \le i \le m$, as $X_i = \langle x_1, x_2, \dots, x_i \rangle$. The same notation is used for other sequences, e.g. Y_j .

Theorem

Let
$$X = \langle x_1, x_2, \dots, x_m \rangle$$
 and $Y = \langle y_1, y_2, \dots, y_n \rangle$ be sequences.
Suppose that $Z = \langle z_1, z_2, \dots, z_k \rangle$ is an LCS of X and Y.

- 1. If $x_m = y_n$, then $z_k = x_m = y_n$ and Z_{k-1} is an LCS for X_{m-1} and Y_{n-1} .
- 2. if $x_m \neq y_n$ and $z_k \neq x_m$, then Z is an LCS for X_{m-1} and Y. 3. If $x_m = y_n$ and $z_k \neq y_n$, then Z is an LCS for X and Y_{n-1} .

Optimal Substructure Theorem

Notation

If $X = \langle x_1, x_2, \dots, x_m \rangle$ is a sequence, then we define X_i , $1 \le i \le m$, as $X_i = \langle x_1, x_2, \dots, x_i \rangle$. The same notation is used for other sequences, e.g. Y_j .

Theorem

Let
$$X = \langle x_1, x_2, \dots, x_m \rangle$$
 and $Y = \langle y_1, y_2, \dots, y_n \rangle$ be sequences.
Suppose that $Z = \langle z_1, z_2, \dots, z_k \rangle$ is an LCS of X and Y.

1. If $x_m = y_n$, then $z_k = x_m = y_n$ and Z_{k-1} is an LCS for X_{m-1} and Y_{n-1} .

2. if
$$x_m \neq y_n$$
 and $z_k \neq x_m$, then Z is an LCS for X_{m-1} and Y

3. If $x_m = y_n$ and $z_k \neq y_n$, then Z is an LCS for X and Y_{n-1} .

< □ ▶ < □ ▶ < □ ▶ < □ ▶ = □ ● ○ < ○

Optimal Substructure Theorem

Notation

If $X = \langle x_1, x_2, \dots, x_m \rangle$ is a sequence, then we define X_i , $1 \le i \le m$, as $X_i = \langle x_1, x_2, \dots, x_i \rangle$. The same notation is used for other sequences, e.g. Y_j .

Theorem

Let
$$X = \langle x_1, x_2, \dots, x_m \rangle$$
 and $Y = \langle y_1, y_2, \dots, y_n \rangle$ be sequences.
Suppose that $Z = \langle z_1, z_2, \dots, z_k \rangle$ is an LCS of X and Y.

1. If $x_m = y_n$, then $z_k = x_m = y_n$ and Z_{k-1} is an LCS for X_{m-1} and Y_{n-1} .

2. if $x_m \neq y_n$ and $z_k \neq x_m$, then Z is an LCS for X_{m-1} and Y. 3. If $x_m \neq y_n$ and $z_k \neq y_n$, then Z is an LCS for X and y_{n-1} .

Proof of the Theorem

Turn this into a solution...

- If x_m = y_n, solve the subproblem
 LCS of X_(m-1), Y_(n-1) and append x_m to get LCS of X, Y.
- If x_m != y_n, solve the subproblems and use the larger solution
 LCS of X_(m-1), Y
 LCS of X, Y_(n-1)

Let C[i,j] be an m-by-n array. C[i,j] will hold the LCS length for X_i, Y_j, where $0 \le i \le m$ and $0 \le j \le n$.

$$C[i_{j},i_{j}] = \begin{cases} 0 & i=0 \text{ or } i=0 \\ c[i_{-1},j-1]+1 & i>0, i>0, x_{i}=y_{j} \\ max\{c[i_{-1},i_{j}], c[i_{j},j-1]\} \\ i>0, x_{i}\neq y_{j} \end{cases}$$



Update min table. Each update is O(i). Solution is O(min) C, G, C, G AnLCS is < G, C, G, C ?.