# DAA - LONGEST COMMON SUBSEQUENCE

https://www.tutorialspoint.com/design\_and\_analysis\_of\_algorithms/design\_and\_analysis\_of\_algorithms\_longest\_common\_subsequence.htm

Copyright © tutorialspoint.com

#### Advertisements

The longest common subsequence problem is finding the longest sequence which exists in both the given strings.

### Subsequence

Let us consider a sequence  $S = \langle s_1, s_2, s_3, s_4, ..., s_n \rangle$ .

A sequence  $Z = \langle z_1, z_2, z_3, z_4, ..., z_m \rangle$  over S is called a subsequence of S, if and only if it can be derived from S deletion of some elements.

### Common Subsequence

Suppose, X and Y are two sequences over a finite set of elements. We can say that Z is a common subsequence of X and Y, if Z is a subsequence of both X and Y.

### Longest Common Subsequence

If a set of sequences are given, the longest common subsequence problem is to find a common subsequence of all the sequences that is of maximal length.

The longest common subsequence problem is a classic computer science problem, the basis of data comparison programs such as the diff-utility, and has applications in bioinformatics. It is also widely used by revision control systems, such as SVN and Git, for reconciling multiple changes made to a revision-controlled collection of files.

#### Naïve Method

Let X be a sequence of length m and Y a sequence of length n. Check for every subsequence of X whether it is a subsequence of Y, and return the longest common subsequence found.

There are  $2^m$  subsequences of X. Testing sequences whether or not it is a subsequence of Y takes On time. Thus, the naïve algorithm would take  $O(n2^m)$  time.

## Dynamic Programming

Let  $X = \langle x_1, x_2, x_3, ..., x_m \rangle$  and  $Y = \langle y_1, y_2, y_3, ..., y_n \rangle$  be the sequences. To compute the length of an element the following algorithm is used.

In this procedure, table C[m, n] is computed in row major order and another table B[m,n] is computed to construct optimal solution.

```
Algorithm: LCS-Length-Table-Formulation (X, Y)
m := length(X)
n := length(Y)
for i = 1 to m do
    C[i, 0] := 0
for j = 1 to n do
```

1 of 2 3/28/2019, 10:31 PM

```
 \begin{array}{l} \texttt{C[0, j]} := 0 \\ \texttt{for } \texttt{i} = 1 \texttt{ to m } \texttt{ do} \\ \texttt{for } \texttt{j} = 1 \texttt{ to n } \texttt{ do} \\ \texttt{if } \texttt{x_i} = \texttt{y_j} \\ \texttt{C[i, j]} := \texttt{C[i - 1, j - 1]} + 1 \\ \texttt{B[i, j]} := \texttt{`D'} \\ \texttt{else} \\ \texttt{if } \texttt{C[i - 1, j]} \ge \texttt{C[i, j - 1]} \\ \texttt{C[i, j]} := \texttt{C[i - 1, j]} + 1 \\ \texttt{B[i, j]} := \texttt{`U'} \\ \texttt{else} \\ \texttt{C[i, j]} := \texttt{C[i, j - 1]} \\ \texttt{B[i, j]} := \texttt{`L'} \\ \texttt{return C and B} \\ \end{array}
```

```
Algorithm: Print-LCS (B, X, i, j)
if i = 0 and j = 0
   return
if B[i, j] = 'D'
   Print-LCS(B, X, i-1, j-1)
   Print(x<sub>i</sub>)
else if B[i, j] = 'U'
   Print-LCS(B, X, i-1, j)
else
   Print-LCS(B, X, i, j-1)
```

This algorithm will print the longest common subsequence of **X** and **Y**.

## Analysis

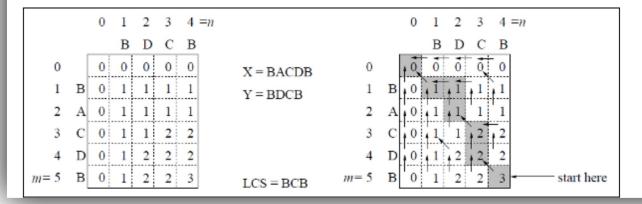
To populate the table, the outer **for** loop iterates m times and the inner **for** loop iterates n times. Hence, the complexity of the algorithm is Om, n, where m and n are the length of two strings.

# Example

In this example, we have two strings X = BACDB and Y = BDCB to find the longest common subsequence.

Following the algorithm LCS-Length-Table-Formulation asstated above, we have calculated table C shown on the lefth and side and table B shown on the right hand side.

In table B, instead of 'D', 'L' and 'U', we are using the diagonal arrow, left arrow and up arrow, respectively. After generating table B, the LCS is determined by function LCS-Print. The result is BCB.



2 of 2 3/28/2019, 10:31 PM