The Longest Common Subsequence problem

K. Subramani LCSEE, West Virginia University, Morgantown, WV {ksmani@csee.wvu.edu}

1 Problem Definition

Let $\mathbf{X} = \langle x_1, x_2, \dots, x_m \rangle$ and $\mathbf{Y} = \langle y_1, y_2, \dots, y_n \rangle$ denote two sequences on a fixed alphabet Σ (You may assume that $\Sigma = \{0, 1\}$, if it helps). Devise an efficient algorithm to determine the length of the Longest Common Subsequence between \mathbf{X} and \mathbf{Y} . Note that a subsequence of a sequence need not be contiguous.

2 Solution

Let $\mathbf{Z} = \langle z_1, z_2, \dots, z_r \rangle$ denote the Longest Common Subsequence (LCS) between \mathbf{X} and \mathbf{Y} . We use the following notational scheme:

- (a) $\mathbf{X}_{\mathbf{i}}$ denotes the sequence $\langle x_1, x_2, \ldots, x_i \rangle$.
- (b) LC(X, Y) denotes the LCS between X and Y.
- (c) f[i][j] denotes the length of the LCS between X_i and Y_j , i.e., $f[i][j] = |LC(X_i, Y_j)|$.

Consider the following cases:

- (i) x_m = y_n In this case, we must have z_r = x_m and further, Z_{r-1} = LC(X_{m-1}, Y_{n-1}). Observe that if z_r ≠ x_m, y_n, then the length of the LCS can be increased by at least 1, by appending x_m to Z, thereby contradicting the optimality of Z. Given that z_r = x_m, it is easy to see that Z_{r-1} = LC(X_{m-1}, Y_{n-1}). If this were not the case, then we can determine LC(X_{m-1}, Y_{n-1}) and append x_m to it, thereby contradicting the optimality of Z.
- (ii) $x_m \neq y_n$ and $z_r \neq x_m$. We must have $\mathbf{Z} = \mathbf{LC}(X_{m-1}, Y_n)$. If this were not the case, then there is a subsequence between \mathbf{X} and \mathbf{Y} of length greater than k. However, this subsequence is also a subsequence between \mathbf{X}_{m-1} and Y, thereby contradicting the optimality of \mathbf{Z} .
- (iii) $x_m \neq y_n$ and $z_r \neq y_n$. We must have $\mathbf{Z} = \mathbf{LC}(X_m, Y_{n-1})$. Symmetric to the case above.

From the above discussion, the following recurrence relation suggests itself for the length of the LCS between X_i and Y_j :

$$\begin{array}{ll} f[i][j] &=& 0, \text{ if } i=0 \text{ or } j=0 \\ &=& \max(f[i-1][j], f[i][j-1]), \text{ if } \mathbf{x_i} \neq \mathbf{y_j} \text{ and } \mathbf{i, j} \geq 1 \\ &=& 1+f[i-1][j-1], \text{ if } \mathbf{x_i} = \mathbf{y_i} \text{ and } \mathbf{i, j} \geq 1 \end{array}$$

This recurrence relation can be easily converted into an algorithm that computes f[m][n] in bottom-up fashion and runs in time $O(m \cdot n)$.