

DNA and Text sequence Alignment

• Sub sequence:

Given a string $X = x_0 x_1 x_2 \dots x_{n-1}$

A subsequence of X is a string S such that

$$(1) \quad S = x_{i_1} x_{i_2} \dots x_{i_k}$$

where $i_j < i_{j+1}$

- not necessarily contiguous
- but must be in order given by X .

• Example sub sequences.

$X = \text{CGATAATTGAGA}$

$S_1 = \text{CGAT}$

$S_2 = \text{ATTT}$

$S_3 = \text{CATT}$

Not subsequence: ACGA
 CGCA

The longest common subsequence (LCS) problem

Problem description:

Given 2 strings:

$$X = x_0 x_1 x_2 \dots x_{n-1}$$

$$Y = y_0 y_1 y_2 \dots y_{m-1}$$

Find the longest ~~sub~~string S such that:

S is a subsequence of X

S is a subsequence of Y

Example :

X = ABCABCABC

Y = BABACBAB

Common subsequence:

AAA

A B C A B C A B C

B A B A C B A B

ABAA

A B C A B C A B C

B A B A C B A B

ABABA

A B C A B C A B C

B A B A C B A B

What is the longest sequence ?

Naive Solution: brute force.

(1) enumerate all subsequences of X

(2) Take the longest one that is a subsequence of Y .

Problem: exponential number of subsequence of X :

$$X = x_0 x_1 x_2 x_3 \dots x_{n-1}$$

↓
step
select or not.

$$2 \cdot 2 \cdot 2 \cdot \dots \cdot 2$$

2^n possible sub sequence!

To test this many will take ages !!!

Eg: $2^{64} = 1844 \dots$

takes 584 years to finish
if one sequence take 10^{-9} sec!

Dynamic Programming (recursion)

If: problem(n) is decomposable into a number of smaller problems:

problem(n_1), problem(n_2), ...

with: $n_1 < n$, $n_2 < n$, ...

And: The solution of problem(n_1), problem(n_2), ... can be used to solve

Problem(n)

EASILY

Then we can apply "Dynamic Programming" (recursion) to solve the problem.

Examples:

$$\text{Fib}(n) = \text{Fib}(n-1) + \text{Fib}(n-2)$$

Hanoi(n) = Hanoi($n-1$) followed by another Hanoi($n-1$)

A dynamic Programming Solution to LCS

• Define: $X[0..i] = x_0 x_1 x_2 \dots x_i$
 $Y[0..j] = y_0 y_1 y_2 \dots y_j$

• Define:

$L[i, j]$ = length of the longest string that
is a subsequence of
 $X[0..i] = x_0 x_1 \dots x_i$
and $Y[0..j] = y_0 y_1 \dots y_j$.

• We want to find $L[n, m]$.

• Define $L[i, j]$ as smaller problems (LCS)
↓
Decompose.

$$X[0..i] = x_0 x_1 \dots x_{i-1} x_i$$

$$Y[0..j] = y_0 y_1 \dots y_{j-1} y_j$$

longest subseq = $L[i, j]$

① $x_i = y_j$:

$$\begin{array}{ccccccc} x_0 & x_1 & \dots & x_{i-1} & x_i \\ y_0 & y_1 & \dots & y_{j-1} & y_j \end{array}$$

match.

longest subsequence in here + $\begin{pmatrix} x_i \\ y_j \end{pmatrix}$

longest subseq.'s length = $L[i-1, j-1]$.

$$\Rightarrow L[i, j] = L[i-1, j-1] + 1.$$

② $x_i \neq y_j$:

$$\begin{array}{ccccccc} x_0 & x_1 & \dots & x_{i-1} & x_i \\ y_0 & y_1 & \dots & y_{j-1} & y_j \end{array}$$

different.

Then we must find $L[i, j]$ matching letters in:

$$\begin{cases} x_0 x_1 \dots x_{i-1} \\ y_0 y_1 \dots y_{j-1} \end{cases} \quad \text{or} \quad \begin{cases} x_0 x_1 \dots x_{i-1} x_i \\ y_0 y_1 \dots y_{j-1} \end{cases}$$

Therefore:

$$L[i,j] = \begin{cases} L[i-1,j] \\ L[i,j-1] \end{cases} \quad \begin{array}{l} \setminus \text{ which ever is} \\ / \text{ larger!} \end{array}$$

$$\Leftrightarrow L[i,j] = \max(L[i-1,j], L[i,j-1]).$$