

[CSE3081(2반)] 알고리즘 설계와 분석

2020학년도 2학기

강의자료

(2020.10.29 목요일)

서강대학교 공과대학 컴퓨터공학과

임 인 성 교수

본 강의에서 제작하여 제공하는 **PDF 파일, 동영상, 그리고 예제 코드 등의 강의 자료**의 저작권은 특별히 명기되어 있지 않은 한 서강대학교에 있습니다.

본인의 학습 목적 외에 공개된 장소에 올리거나 타인에게 배포하는 등의 행위를 금합니다. 협조 부탁드립니다.

[주제 4]

Dynamic Programming

Longest Common Subsequence (LCS)

- Definitions

- Given a sequence $X = \langle x_1, x_2, \dots, x_m \rangle$, another sequence $Z = \langle z_1, z_2, \dots, z_k \rangle$ is a **subsequence** of X if there exists a strictly increasing sequence $\langle i_1, i_2, \dots, i_k \rangle$ of indices of X such that for all $j = 1, 2, \dots, k$, we have $x_{i_j} = z_j$.
 - A subsequence of a given sequence is just the given sequence with some elements (possibly none) left out.
 - Ex: $X = \langle A, B, C, B, D, A, B \rangle$, $Z = \langle B, C, D, B \rangle$ ($\langle 2, 3, 5, 7 \rangle$)
- Given two sequences X and Y , we say that a sequence Z is a **common subsequence** of X and Y if Z is a subsequence of both X and Y .
 - Ex: $X = \langle A, B, C, B, D, A, B \rangle$, $Y = \langle B, D, C, A, B, A \rangle$, $Z_1 = \langle B, C, A \rangle$, $Z_2 = \langle B, C, B, A \rangle$, $Z_3 = \langle B, D, A, B \rangle$
- Given a sequence $X = \langle x_1, x_2, \dots, x_m \rangle$, $X_i = \langle x_1, x_2, \dots, x_i \rangle$ is the i th **prefix** of X , for $i = 0, 1, \dots, m$.
 - Ex: $X = \langle A, B, C, B, D, A, B \rangle$, $X_4 = \langle A, B, C, B \rangle$, $X_0 = \text{null sequence}$

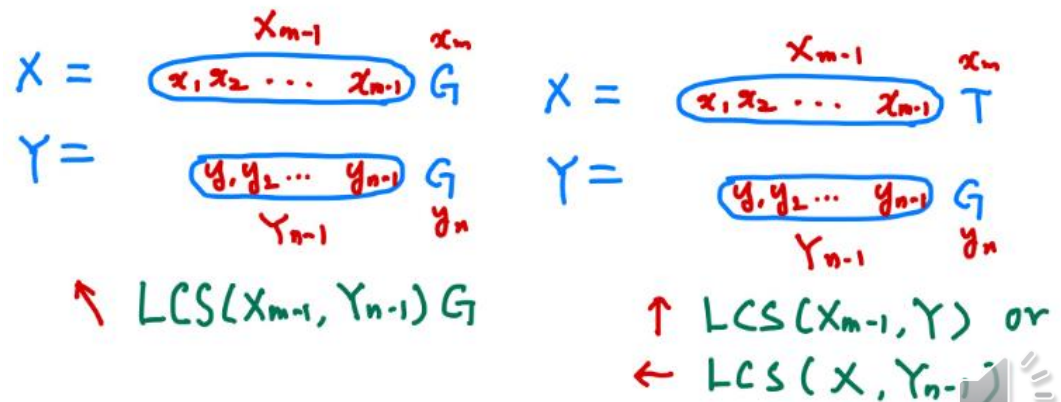
• Problem

- Given two sequences $X = \langle x_1, x_2, \dots, x_m \rangle$ and $Y = \langle y_1, y_2, \dots, y_n \rangle$, find a *longest common subsequence* of X and Y .

X = TCCCCGCTCTGCTCTGTCCGGTCACAGGACTTTTTGCCCTCTGTTCCCGGGTCCCTCAGGCGGCCACCCA
GTGGGCACACTCCCAGGCGGCGCTCCGGCCCCGCGCTCCCTCCCTCTGCCTTTTATTCCCAGCTGTCAAC
ATCCTGGAAGCTTTGAAGCTCAGGAAAGAAGAGAAATCCACTGAGAACAGTCTGTAAAGGTCCGTAGTGC
TATCTACATCCAGACGGTGGAAGGGAGAGAAAGAGAAAGAGGTATCCTAGGAATACCTGCCTGCTTAGA
CCCTCTATAAAAGCTCTGTGCATCCTGCCACTGAGGACTCCGAAGAGGTAGCAGTCTTCTGAAAGACTTC
AACTGTGAGGACATGTCGTTTCAAGTTTGGCCAACATCTCATCAAGCCCTCTGTAGTGTCTTCTCAAAACAG
AACTGTCCTTCGCTCTTGTGAATAGGAAACCTGTGGTACCAGGACATGTCTTGTGTGCCCGCTGCGGCC
AGTGGAGCGCTTCCATGACCTGCGTCTGATGAAGTGGCCGATTTGTTTCAGACGACCCAGAGAGTCCGGC
ACAGTGGTGGAAAAACATTTCCATGGGACCTCTCTCACCTTTTCCATGCAGGATGGCCCCGAAGCCGGAC
AGACTGTGAAGCACGTTTACGTCCTATGTTCTTCCAGGAAGGCTGGAGACTTTTACAGGAATGACAGCAT
CTATGAGGAGCTCCAGAAACATGACAAGGAGGACTTTTCTGCCTCTTGGAGATCAGAGGAGGAAATGGCA
GCAGAAGCCGCAGCTCTGCGGGTCTACTTTTCAAGTACACAGATGTTTTTTCAGATCCTGAATTCCAGCAA
AGAGCTATTGCCAACAGTTTGAAGACCGCCCCCGCCTCTCCCCAAGAGGAACTGAATCAGCATGAAA
ATGCAGTTTCTTATCTCACCATCCTGTATTCTTCAACAGTGTATCCCCACCTCGGTCACTCCAACTCC
CTTAAATACCTAGACCTAAACGGCTCAGACAGGCAGATTTGAGGTTTCCCCCTGTCTCCTTATTCGGCA
GCCTTATGATTAACTTCCTTCTCTGCTGCAAAAAAAAAAAAAAAAAA

Y = ATGTTAACCAAGGAATGGATCTGTGTCGTTCCAGTTTGAAGGCCTTTTCTGATGAAATGAAGATAGGTT
TCAACTCCACAGGTTATTGTGGTATGATCTTAACCAAAAATGATGAAGTTTTCTCCAAGATTACTGAAAA
ACCTGAATTGATTAACGATATCTTATTGGAATGTGGTTTCCCAAACACTTCTGGTCAAAAACCAACGAA
TACAACTATTGAGTCCTTACAAGTACAAGTATACACGTATACATGTATGTATATATATATGATTTAAA
TGATAGAAGTAATTTCTATATGTATATGTCTATTCAATTTTTATTCTAATGACTTTGAAATTTTATATT
TATTATTCTACACTTTATTATTAATAACTACATGCAATCAATGCCGCTAAGGTTACGATTCACCTTTTTAT
TACTTATTATTAATATTGTTGTATTACTTCTTGAATAATATGTCTAAAGAGTCCTAATTTGGATTTTC
TTTTCTCCTAACTTCACTGTCTGCGCTGCTTTTCTAACGCACCATCGCTAATACACCAGCTTTTATT
GCTTGTGCTGCGCTATTGCTCGCATGGAACGTTTTTCAAGTGCCTCATCATCCTGGGATAAACTAAAGACT
AAGTCACCAGTTTCAATTTGAGGCTTTTCTCTGCGTGTGACAAAAGAGGACAGATCAACCATATCACCTG
GTTACCATGAGAATGTCCTTACTTAGTTGCGAATTTGGTTCTGATCTGCCTTCAGACTTGAATCATTC
ATTACTGTCATTACTTTTAGCCTATTCATTTTCTTCTTGCCTATCAGGTACAGGGATTTGACCACAGAGT
GTTGAAGGGGTGCATCGTCCGCTCTCGTAATAACCCTCCGATACTATTTTCAATGTTGGCACGTTGCACT
GAAAAGGGCACTTGGCACTGTGCACTTTTAATGTTTTTCAATTTTTCATCATATCATCATATGGCATTTCAT
AATGTTGACTCTTGTAGTTGAAGATTGAGTTTATCGTGTTACTGTTTCTGCTACCTTTTCAATTATCAAT
CAGGTTGCGGTGTTGCATGTGGGAGATAGGACTGCCGATCTTCTTCTTCTCGATTCTCCACTAAATCC
TGCTTTTTATCGCTATCCAGCACACGATTGAGCTGTGAATTGCCGCACTTTTTAGGATAACCATCCTTGG

- Naïve approach
 - Enumerate all subsequences of X and check each subsequence to see if it is also a subsequence of Y , keeping track of the longest subsequence found.
 - *Exponential algorithm!*
 - ✓ The LCS problem can be solved efficiently using dynamic programming.
- **Optimal substructure of an LCS:** Let $X = \langle x_1, x_2, \dots, x_m \rangle$ and $Y = \langle y_1, y_2, \dots, y_n \rangle$ be sequences, and let $Z = \langle z_1, z_2, \dots, z_k \rangle$ be any 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 of X_{m-1} and Y_{n-1} .
 - ② If $x_m \neq y_n$, then an LCS of X and Y is either an LCS of X_{m-1} and Y or an LCS of X and Y_{n-1} .



- Let $c[i, j]$ be the length of an LCS of the sequences X_i and Y_j .
- Optimal substructure for computing $c[i, j]$

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

$$\begin{array}{l} X = \overset{x_{m-1}}{\underbrace{x_1 x_2 \dots x_{n-1}}} \overset{x_m}{G} \\ Y = \underbrace{y_1 y_2 \dots y_{n-1}}_{Y_{n-1}} \overset{y_n}{G} \\ \quad \uparrow \text{LCS}(X_{m-1}, Y_{n-1}) \end{array}$$

$$\begin{array}{l} X = \overset{x_{m-1}}{\underbrace{x_1 x_2 \dots x_{n-1}}} \overset{x_m}{T} \\ Y = \underbrace{y_1 y_2 \dots y_{n-1}}_{Y_{n-1}} \overset{y_n}{G} \\ \quad \uparrow \text{LCS}(X_{m-1}, Y) \text{ or} \\ \quad \leftarrow \text{LCS}(X, Y_{n-1}) \end{array}$$

$O(mn)$ Algorithm

- Filling the table

$$X = \overbrace{x_1 x_2 \dots x_{n-1}}^{x_{n-1}} x_n G$$

$$Y = \overbrace{y_1 y_2 \dots y_{n-1}}^{y_{n-1}} y_n G$$

↑ $LCS(X_{n-1}, Y_{n-1}) G$

$$X = \overbrace{x_1 x_2 \dots x_{n-1}}^{x_{n-1}} x_n T$$

$$Y = \overbrace{y_1 y_2 \dots y_{n-1}}^{y_{n-1}} y_n G$$

↑ $LCS(X_{n-1}, Y)$ or
← $LCS(X, Y_{n-1})$

LCS-LENGTH(X, Y)

```

1   $m \leftarrow \text{length}[X]$ 
2   $n \leftarrow \text{length}[Y]$ 
3  for  $i \leftarrow 1$  to  $m$ 
4      do  $c[i, 0] \leftarrow 0$ 
5  for  $j \leftarrow 0$  to  $n$ 
6      do  $c[0, j] \leftarrow 0$ 
7  for  $i \leftarrow 1$  to  $m$ 
8      do for  $j \leftarrow 1$  to  $n$ 
9          do if  $x_i = y_j$ 
10             then  $c[i, j] \leftarrow c[i - 1, j - 1] + 1$ 
11                   $b[i, j] \leftarrow \text{"\textbackslash"}$ 
12             else if  $c[i - 1, j] \geq c[i, j - 1]$ 
13                 then  $c[i, j] \leftarrow c[i - 1, j]$ 
14                         $b[i, j] \leftarrow \text{"\textup{\textsf{↑"}}$ 
15                 else  $c[i, j] \leftarrow c[i, j - 1]$ 
16                         $b[i, j] \leftarrow \text{"\textleftarrow"}$ 
17  return  $c$  and  $b$ 
    
```

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

j		0	1	2	3	4	5	6
i	y_j	B	D	C	A	B	A	
	x_i							
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

- Printing the LCS

$X = \overset{x_{m-1}}{\underbrace{x_1, x_2, \dots, x_{m-1}}} \overset{x_m}{G}$
 $Y = \underbrace{y_1, y_2, \dots, y_{n-1}}_{Y_{n-1}} \overset{y_n}{G}$
 $\uparrow \text{LCS}(X_{m-1}, Y_{n-1}) G$

$X = \overset{x_{m-1}}{\underbrace{x_1, x_2, \dots, x_{m-1}}} \overset{x_m}{T}$
 $Y = \underbrace{y_1, y_2, \dots, y_{n-1}}_{Y_{n-1}} \overset{y_n}{G}$
 $\uparrow \text{LCS}(X_{m-1}, Y) \text{ or } \leftarrow \text{LCS}(X, Y_{n-1})$

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

```

1  if  $i = 0$  or  $j = 0$ 
2      then return
3  if  $b[i, j] = "\diagup"$ 
4      then PRINT-LCS( $b, X, i - 1, j - 1$ )
5          print  $x_i$ 
6  elseif  $b[i, j] = "\uparrow"$ 
7      then PRINT-LCS( $b, X, i - 1, j$ )
8  else PRINT-LCS( $b, X, i, j - 1$ )
    
```

j		0	1	2	3	4	5	6
i	y_j		B	D	C	A	B	A
	x_i							
0		0	0	0	0	0	0	0
1	A	0	\uparrow 0	\uparrow 0	\uparrow 0	\swarrow 1	\leftarrow 1	\swarrow 1
2	B	0	\swarrow 1	\leftarrow 1	\leftarrow 1	\uparrow 1	\swarrow 2	\leftarrow 2
3	C	0	\uparrow 1	\uparrow 1	\swarrow 2	\leftarrow 2	\uparrow 2	\uparrow 2
4	B	0	\swarrow 1	\uparrow 1	\uparrow 2	\uparrow 2	\swarrow 3	\leftarrow 3
5	D	0	\uparrow 1	\swarrow 2	\uparrow 2	\uparrow 2	\uparrow 3	\uparrow 3
6	A	0	\uparrow 1	\uparrow 2	\uparrow 2	\swarrow 3	\uparrow 3	\swarrow 4
7	B	0	\swarrow 1	\uparrow 2	\uparrow 2	\uparrow 3	\swarrow 4	\uparrow 4

C Implementation

Courtesy of
<http://www.bioalgorithms.info/downloads/code/>

```
/** Copyright (C) 2005 Neil Jones. */

#include <stdio.h>
char* LCS(char* a, char* b);

#define NEITHER 0
#define UP 1
#define LEFT 2
#define UP_AND_LEFT 3

int main(int argc, char* argv[]) {
    printf("%s\n", LCS(argv[1],argv[2]));
}

char* LCS(char* a, char* b) {
    int n = strlen(a); int m = strlen(b);
    int** S; int** R; int ii; int jj;
    int pos; char* lcs;

    S = (int **)malloc( (n+1) * sizeof(int *) );
    R = (int **)malloc( (n+1) * sizeof(int *) );

    for (ii = 0; ii <= n; ++ii) {
        S[ii] = (int*) malloc( (m+1) * sizeof(int) );
        R[ii] = (int*) malloc( (m+1) * sizeof(int) );
    }
```

```
    for (ii = 0; ii <= n; ++ii) {
        S[ii][0] = 0; R[ii][0] = UP;
    }
    for (jj = 0; jj <= m; ++jj) {
        S[0][jj] = 0; R[0][jj] = LEFT;
    }
    for (ii = 1; ii <= n; ++ii) {
        for (jj = 1; jj <= m; ++jj) {
            if ( a[ii-1] == b[jj-1] ) {
                S[ii][jj] = S[ii-1][jj-1] + 1;
                R[ii][jj] = UP_AND_LEFT;
            }
            else {
                S[ii][jj] = S[ii-1][jj-1] + 0;
                R[ii][jj] = NEITHER;
            }
            if ( S[ii-1][jj] >= S[ii][jj] ) {
                S[ii][jj] = S[ii-1][jj];
                R[ii][jj] = UP;
            }
            if ( S[ii][jj-1] >= S[ii][jj] ) {
                S[ii][jj] = S[ii][jj-1];
                R[ii][jj] = LEFT;
            }
        }
    }
}
```

```

ii = n;
jj = m;
pos = S[ii][jj];
lcs = (char *) malloc( (pos+1) * sizeof(char) );
lcs[pos--] = (char)NULL;

while ( ii > 0 || jj > 0 ) {
    if( R[ii][jj] == UP_AND_LEFT ) {
        ii--; jj--;
        lcs[pos--] = a[ii];
    }
    else if ( R[ii][jj] == UP ) {
        ii--;
    }
    else if ( R[ii][jj] == LEFT ) {
        jj--;
    }
}
for (ii = 0; ii <= n; ++ii ) {
    free(S[ii]); free(R[ii]);
}
free(S);
free(R);
return lcs;
}

```

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

The Gapped Alignment Problem

$A = \varepsilon$
 $B = ATCG$

$\overline{\overline{A}} \quad \overline{\overline{T}} \quad \overline{\overline{C}} \quad \overline{\overline{G}}$

- Problem:** Given two sequences, find a gapped alignment that maximize the score!

- Compare two sequences if they are similar (related).
- Gapped alignment

- Example: $A = \text{ATCGGATCT}$, $B = \text{ACGGACT}$

A	T	C	G	G	A	T	-	C	T	A	T	C	G	G	A	T	C	T
A	-	C	-	G	G	-	A	C	T	A	-	C	G	G	-	A	C	T

$$5 * 2 + 1 * (-1) + 4 * (-2) = 1 \quad 6 * 2 + 1 * (-1) + 2 * (-2) = 7$$

- A possible alignment scoring scheme
 - Ex: match score = 2, mismatch penalty = -1, gap penalty = -2

Optimal substructure

$$S(i, j) = \begin{cases} -2j, & \text{if } i = 0 \\ -2i, & \text{if } j = 0 \\ \max\{ S(i-1, j-1) + s(a_i, b_j), S(i, j-1) - 2, S(i-1, j) - 2 \}, & \text{if } i > 0 \text{ and } j > 0 \end{cases}$$

$A_i = a_1 a_2 \dots a_{i-1} a_i$

$B_j = b_1 b_2 \dots b_{j-1} b_j$

$A_i = a_1 a_2 \dots a_{i-1} a_i$

$B_j = b_1 b_2 \dots b_{j-1} b_j$

$A_i = a_1 a_2 \dots a_{i-1} a_i$

$B_j = b_1 b_2 \dots b_{j-1} b_j$

Longest Increasing Subsequence (LIS)

- **Problem:** Given a sequence $A = (a[0], a[1], \dots, a[n-1])$, find the length of the longest subsequence such that all elements of the subsequence are sorted increasing order.
- **Example**
 - $(10, 22, 9, 33, 21, 50, 41, 60, 80) \rightarrow (10, 22, 33, 50, 60, 80)$
 - $(0, 8, 4, 12, 2, 10, 6, 14, 1, 9, 5, 13, 3, 11, 7, 15) \rightarrow (0, 2, 6, 9, 11, 15)$
 $(0, 4, 6, 9, 11, 15)$

- **Algorithm**

Let $d[i]$ be the length of the longest increasing subsequence that ends in the element at index i . Then, the answer to the LIS problem is the maximum value of $d[i]$, $i = 0, 1, \dots, n - 1$.

- **Optimal substructure**

$$d[i] = \max \left(1, \max_{\substack{j=0 \dots i-1 \\ a[j] < a[i]}} (d[j] + 1) \right), i = 0, 1, \dots, n-1$$

(0, 8, 4, 12, 2, 10, 6, 14, 1, 9
 \downarrow
 $ac[i]$

```

int LIS( int* a, int N ) {
    int *best, *prev, i, j, max = 0;

    best = (int*) malloc ( sizeof( int ) * N );
    prev = (int*) malloc ( sizeof( int ) * N );
    for ( i = 0; i < N; i++ ) best[i] = 1, prev[i] = i;

    for ( i = 1; i < N; i++ )
        for ( j = 0; j < i; j++ )
            if ( a[i] > a[j] && best[i] < best[j] + 1 )
                best[i] = best[j] + 1, prev[i] = j;

    for ( i = 0; i < N; i++ )
        if ( max < best[i] ) max = best[i];
    // Print the LIS using prev[] here.
    free( best ); free( prev );
    return max;
}

```

$$(\underline{0}, \underline{8}, \underline{4}, \underline{12}, \underline{2}, \underline{10}, \underline{6}, \underline{14}, \underline{1}, \underline{9})$$

$$\downarrow$$

$$a[i]$$

$$d[i] = \max \left(1, \max_{\substack{j=0 \dots i-1 \\ a[j] < a[i]}} (d[j] + 1) \right), i = 0, 1, \dots, n-1$$