6.6 Longest Common Subsequence

A strand of DNA consists of a string of molecules called bases, which can be one of adenine, cytosine, guanine, and thymine. If we use the first letter to represent each base, we can express a strand of DNA as a string over the 4-element set {A, C, G, T}. We might have two organisms whose DNA strands are

$S_1 = ACTGCGTCAGCTGCAGCGTCGA$ $S_2 = CAGCTGATCGTAGCTGATCG$

There are a few ways one may decide how similar two strings like the above are. One way is to find the longest common substring. Another way is to count the number of characters we have to edit from string 1 to get string 2, also known as the edit distance. There is a DP solution to find the edit distance, but we will not focus on this solution today.

What we will attempt to find is the longest common subsequence between the two strings. Don't confuse this with the longest common substring! In the longest common subsequence, the goal is to find a third string whose characters appear in the same order in both S_1 and S_2 .

Question 96. Verify that ACTGCGTAGCTGATCG is a valid subsequence of *S*₁ and *S*₂.

$S_{1} = ACTGCGTCAGCTGCAGCGTCGA$ ACTGCGTAGCTGATCG $S_{2} = CAGCTGATCGTAGCTGATCG$

Problem (LCS).

INPUT: Two sequences *X* with length *m* and *Y* with length *n*.

OUTPUT: The length of the maximum length common subsequence of *X* and *Y*.

Question 97. Reason about whether or not this problem has the optimal substructure property.

- Write down a concise statement for the "optimal solution" using the input parameters.
- What are a set of easy top level choices to make, which break the problem into multiple subproblems? Remember, a subproblem must also be a valid instance of the main problem!

Question 98. Using the above analysis, write down a recursive algorithm that solves the problem. This algorithm should not use a memo and will likely be inefficient.

Question 99. Given your answer to the previous problem, does the LCS problem also have overlapping subproblems? If yes, write down a scenario where the same subproblem gets called twice.

Question 100. Identify the pieces of the dynamic programming algorithm.

Subproblem domain: The space of possible subproblems to consider.

Memo table definition: The items we will be storing in our tables. Alternatively, the name of solutions to subproblems.

Goal: The location of our solution in the table.

Initial values: Solutions to trivial subproblems to start building on.

Recurrence: A compact representation of the recursive function.

The following is some python code to compute the LCS length in a bottom up fashion.

```
def LCS(X, Y):
1
       m, n = len(X), len(Y)
2
       memo = # Table with [0:m, 0:n]. Stores the LCS of X[1:m], Y[1:n]
3
       prev = # Table with [1:m, 1:n].
4
                # Stores the subproblem this solution builds out of.
5
       for i in range(1, m + 1):
6
           memo[i][0] = 0
7
       for j in range(0, n + 1):
8
           memo[0][j] = 0
9
       for i in range(1, m + 1):
10
            for j in range(1, n + 1):
11
                if X[i] == Y[j]:
12
                    memo[i][j] = memo[i - 1][j - 1] + 1
13
                    prev[i][j] = "\uparrow \leftarrow "
14
                elif memo[i - 1][j] >= memo[i, j - 1]:
15
                    memo[i][j] = memo[i - 1][j] + 1
16
                    prev[i][j] = "^"
17
                else:
18
                    memo[i][j] = memo[i][j - 1] + 1
19
                    prev[i][j] = "←"
20
       return memo, prev
21
```

Question 101. Draw an example that corresponds to each of the branches in the if statement.

The following is an example table that you may have filled after the inputs X = ABCBDAB and Y = BDCABA.

	j	0	1	2	3	4	5	6
i		y _j	В	D	С	А	В	Α
0	x_i							
		0	0	0	0	0	0	0
1	A		1	1	Ť		←	1
		0	0	0	0		1	1
2	В		~	\leftarrow	\leftarrow	1	~	\leftarrow
		0	1	1	1	1	2	2
3	C		1	1	~		1	1
		0	1	1	2		2	2
4	В		~	1	1	1	~	\leftarrow
		0	1	1	2	2	3	3
5	D		1	5	1	1	1	
		0	1	2	2	2	3	
6	A		1	1	1	5	1	5
		0	1	2	2	3	3	4
7	В		5	1	1	1	5	
		0	1	2	2	3	4	

Question 102. Fill out the cells that are empty in the table.

Question 103. Can you use the table above to reconstruct the LCS?