

## Problem Set

Please be kind while grading :)

### Problem 0 – Collaboration Policy (10%)

Read and sign the Collaboration Policy in the Course Content folder on HuskyCT. Hand in your signed form with your homework submission. We cannot grade your work without a signed form.

**You've already responded**

You can fill out this form only once.

Try contacting the owner of the form if you think this is a mistake.

### Problem 1 – Introductions (10%)

The goal of this problem is to introduce yourself to the course staff and to think about what you want to get out of the course. Write at most two paragraphs about yourself. This may include your academic and non-academic interests, preferred programming languages, courses that you are excited to take in the future, or areas of computer science that you find the most interesting. Finish by describing what do you want to learn by taking this course. Think about what would be most beneficial to your prospective career.

#### Introduction:

Hi my name is Benny and I am a Junior majoring in Computer Science. I mainly am taking this class for my major however I know the importance of learning Algorithms and Complexities in Computer Science so I want to learn as much as I can from this class. I currently am a Data Analytics concentration as I like looking at and manipulating data. I somewhat like algorithms however I am not the best at proving them like with complexities associated with them. I personally like programming in python as it is easy to use and versatile enough for smaller tasks, however I also like programming in C++. I am most interested in machine learning in the computer science field as I love the applications and possibilities of

it. From this class my goal is to learn how to prove algorithms and their complexities as well as learn more about the different types of algorithms.

**Problem 2 – Preliminaries (10%)**

This problem tests your retention of prerequisite materials. If some of these problems are difficult, please go back and review material from the prerequisites.

**Questions:**

1. How many permutations are there of the set of  $n$  numbers?
2. How many subsets of three numbers are there from a set of  $n$  numbers?
3. How many subsets of  $S$  are there when  $|S| = n$ ?

**Answer:**

1. Of a set of  $n$  numbers, there are going to be  $n!$  permutations. As there are  $n$  numbers, there are  $n$  ways to choose the first number,  $n - 1$  ways to choose the second number, and so on. Each of these choices is independent of the others, so the total number of permutations is the product of the number of choices for each position.
2. There are going to be  $n^3$  subsets of three numbers from a set of  $n$  numbers.
3. There are going to be  $2^n$  subsets of  $S$  when  $|S| = n$ .

**Problem 3 – Complexity (20%)**

1. For functions  $A$  and  $B$  and constant  $c > 1$ , indicate which of  $\{A = O(B), A = \Omega(B), A = \Theta(B)\}$  holds. Here, we use the notation that  $\log_2 = \lg$ . *Hint: for these problems it is useful to know some mathematical identities. See CLRS section 3.2 or **standard notations and common functions.pdf** in HuskyCT Lecture 2 course notes.*

(a)  $A = n^c, B = c^n$

(b)  $A = n^{\lg(c)}, B = c^{\lg(n)}$

(c)  $A = \lg(n!), B = \lg(n^n)$ . For this problem, you will find Stirling's approximation useful:  $n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \left(1 + \Omega\left(\frac{1}{n}\right)\right)$ .

(d)  $A = 3^{3^n}, B = 3^{n^2}$

2. Place the following functions in increasing order of growth rate. I.e. if  $g$  follows  $f$  then  $f = O(g)$  is true.

(a)  $n^e$ ,

(b)  $n^{\lg(n)}$

(c)  $n^\pi$ ,

(d)  $\lg(n!)$

(e)  $2\sqrt{\lg(n)}$

(f)  $n$

**Answer:**

1. (a)  $A = O(B)$

(b)  $A = \Theta(B)$

(c)  $A = O(B)$

(d)  $A = \Omega(B)$

2. (a)  $2\sqrt{\lg(n)}$

(b)  $n$

(c)  $\lg(n!)$

(d)  $n^e$

(e)  $n^\pi$

(f)  $n^{\lg(n)}$

## Problem 4 – Vinder (30%)

Vinder is a new smartphone app designed to make it easier for everyone to eat their vegetables. Vinder compares your genetic information with a large database of vegetable genetic sequences to match you with the right vegetable based on sequence similarity. Specifically, it computes the longest common subsequence between your DNA sequence and the DNA sequences of thousands of vegetables.

An  $n$ -length *DNA sequence*,  $A$ , is an ordered collection of elements  $\in [A, C, G, T]$ , or  $A \in [A, C, G, T]^n$ . A *subsequence* is a sequence formed by deleting elements from another sequence, thus preserving the original order. A  $k$ -length *prefix* of sequence  $A$  is the first  $k$  elements of  $A$ . For example,  $A = [A, C, C, G, G, A, A, T, C]$  is a sequence.  $[A, C, A, T]$  and  $[C]$  are subsequences of  $A$  but  $[T, A]$  is not.  $[A, C, C, G]$  is a 4-prefix of  $A$  and  $[]$  is the 0-prefix.

Our goal is to find the longest common subsequence (LCS) of  $A$ , the human DNA sequence, and sequences  $B^k \in D$  where  $D$  is the database of vegetable DNA sequences and  $k = 1, \dots, |D|$ . This solution can be found efficiently using dynamic programming. At a high-level, the algorithm is:

1. For each sequence  $B^k \in D$
2. Initialize the length of the LCS to 0.
3. Let  $A_i$  denote the  $i^{\text{th}}$  element of  $A$ . Consider each pair of  $A_i$  and  $B_j^k$ ; either  $A_i = B_j^k$  or  $A_i \neq B_j^k$ .
  - If  $A_i = B_j^k$ , then the length of the LCS of the  $i$ -prefix of  $A$  and  $j$ -prefix of  $B^k$  is one more than the LCS of the  $i - 1$ -prefix of  $A$  and  $j - 1$ -prefix of  $B^k$ .
  - If  $A_i \neq B_j^k$ , then we cannot extend the LCS. Instead, we conclude that the LCS up to  $(A_i, B_j^k)$  is the maximum of the LCS up to  $(A_{i-1}, B_j^k)$  or  $(A_i, B_{j-1}^k)$ .

The two aforementioned cases sum up the possibilities at  $(A_i, B_j^k)$ . We can use these to recursively define the LCS in terms of smaller problem solutions. Your goal is to describe the algorithm to do this.

1. Describe an algorithm to find the length of the LCS that does not use dynamic programming. You do not have to write pseudocode, simply describe the algorithm in enough detail so that we understand it. What is the runtime? Why is it correct? You may keep your answers informal. *Hint: a brute-force or recursive solution should have exponential runtime.*

**Answer:** A algorithm that does not use dynamic programming to solve this problem would be to try all possible subsequences of  $A$  and  $B^k$  and find the longest one. We can do this by checking if the first element of  $A$  is the same as the first element of  $B^k$

and if it is, we can check the rest of the subsequences of  $A$  and  $B^k$  by removing the first element of both sequences. This would have exponential runtime because there are  $2^n$  possible subsequences of  $A$  and  $2^m$  possible subsequences of  $B^k$  where  $n$  and  $m$  are the lengths of  $A$  and  $B^k$  respectively. This would be a brute force solution and would not be efficient.

2. Write out the recurrence relation used in the algorithm defined in the problems description (not part 1 of the answers). That is, write out how the LCS length can be expressed in terms of smaller instances of the problem.

**Answer:** Let  $LCS(i, j)$  be the length of the LCS of the  $i$ -prefix of  $A$  and  $j$ -prefix of  $B^k$ . Then, 
$$LCS(i, j) = \begin{cases} 0 & \text{if } i = 0 \text{ or } j = 0 \\ LCS(i - 1, j - 1) + 1 & \text{if } A_i = B_j^k \\ \max(LCS(i - 1, j), LCS(i, j - 1)) & \text{if } A_i \neq B_j^k \end{cases}$$

3. Think about how to convert this into a dynamic program. What size table would you need for the LCS problem?

**Answer:** We would need a table of size  $n \times m$  where  $n$  is the length of  $A$  and  $m$  is the length of  $B^k$ .

4. Interpret each table entry. That is, give your explanation for what each entry in the dynamic programming table represents. We discussed how to interpret the entries of the dynamic programming table for the similar problem of *edit distance*.

**Answer:** Each entry in the table represents the length of the LCS of the  $i$ -prefix of  $A$  and  $j$ -prefix of  $B^k$ .

5. Write out the table for *Derek* :  $[A, C, A, G, G, T, T, A, C]$  and *Asparagus*<sup>1</sup> :  $[T, C, G, G, A, A, T, A, A]$ .

**Answer:**

	A	C	A	G	G	T	T	A	C
T	0	0	0	0	0	1	1	1	1
C	0	1	1	1	1	1	1	1	2
G	0	1	1	2	2	2	2	2	2
G	0	1	1	2	3	3	3	3	3
A	1	1	2	2	3	3	3	4	4
A	1	1	2	2	3	3	3	4	4
T	1	1	2	2	3	4	4	4	4
A	1	1	2	2	3	4	4	5	5
A	1	1	2	2	3	4	4	5	5

---

<sup>1</sup>provably the best vegetable

6. Prove your LCS algorithm's (that you defined in 2-5) correctness. *Hint: Many dynamic programming proofs follow a simple structure. Consider the following template:*
- (a) Describe what you will show. This should include the variables used to index the table, presented in the order in which your algorithm fills the table values. Also detail here what your algorithm is supposed to compute.
  - (b) State your induction hypothesis. This includes stating your base case and how your algorithm correctly computes the base cases. Then describe an arbitrary input and corresponding arbitrary entry in the table that will be computed. You can think of this as a *state* of the algorithm at one particular time point. Next, consider an optimal dynamic programming table for your particular problem; detail the last decision made in this optimal solution. This usually falls under a number of different cases based on your recurrence relation.
  - (c) Consider each case; assuming that everything in the dynamic programming table is correct, prove that your algorithm will correctly compute the current entry in the table. Show that your recurrence relation accurately describes each possibility for the last choice made by the unknown LCS. This is the part of your proof where you compare the real world possibilities to the cases in your code. By showing this, your recurrence relation and code rely on the previous entries in the table to produce a solution at least as good as the optimal LCS.

**Answer:**

- (a) We will show that the algorithm correctly computes the length of the LCS of  $A$  and  $B^k$ . The algorithm will fill the table in a bottom-up fashion. It will fill the first row and column with 0's. Then, it will fill the rest of the table by considering the case where  $A_i = B_j^k$  and the case where  $A_i \neq B_j^k$ . In the case where  $A_i = B_j^k$ , the algorithm will fill the table entry with the value of the table entry in the previous row and column plus 1.
- (b) Our induction hypothesis is that the algorithm correctly computes the length of the LCS of the  $i$ -prefix of  $A$  and  $j$ -prefix of  $B^k$  for all  $i, j$ .
- (c) The algorithm will correctly compute the length of the LCS of the  $i$ -prefix of  $A$  and  $j$ -prefix of  $B^k$  by taking the maximum of the length of the LCS of the  $(i-1)$ -prefix of  $A$  and  $j$ -prefix of  $B^k$  and the length of the LCS of the  $i$ -prefix of  $A$  and  $(j-1)$ -prefix of  $B^k$ . This is because the algorithm will correctly compute the length of the LCS of the  $(i-1)$ -prefix of  $A$  and  $j$ -prefix of  $B^k$  and the length of the LCS of the  $i$ -prefix of  $A$  and  $(j-1)$ -prefix of  $B^k$  by the induction hypothesis. The base case is that the algorithm correctly computes the length of the LCS of the 0-prefix of  $A$  and 0-prefix of  $B^k$  which is 0. In this case, the algorithm will correctly compute the length of the LCS of the  $i$ -prefix of  $A$  and  $j$ -prefix of  $B^k$  by the induction hypothesis. In the case where  $A_i = B_j^k$ , the algorithm will correctly



compute the length of the LCS of the  $i$ -prefix of  $A$  and  $j$ -prefix of  $B^k$  by taking the value of the table entry in the previous row and column plus 1.

## Problem 5 – Longest Path Problem (20%)

We are given a directed, unweighted graph  $G = (V, E)$  where  $v_i \in V$  for  $i = 1, \dots, n$ . Additionally,  $G$  is ordered, meaning

- An edge can only originate from a vertex with lower index than its destination. Formally, edges may only have the form  $(v_i, v_j)$  where  $i < j$ .
- The only *absorbing* node is  $v_n$ . Absorbing nodes have out-degree 0.

Given an ordered, unweighted, directed graph  $G(V, E)$ , the *longest path problem* is to compute the longest path length from  $v_1$  to  $v_n$ .

Algorithm 1: Longest path algorithm

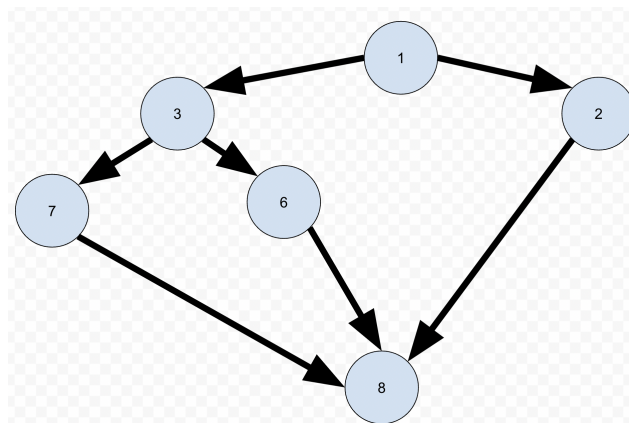
---

```
def longest_path(v1):
    iter=v1
    length=0
    while iter.hasOutgoingEdge():
        iter=iter.getEdgeSmallestIndex()
        length+=1
    return length
```

---

1. Demonstrate that Alg. 2 does not correctly solve the problem by giving a counterexample as a graph. The function `iter.hasOutgoingEdge()` returns a boolean indicating whether or not the vertex `iter` has an outgoing edge. The function `iter.getEdgeSmallestIndex()` returns the smallest index of all vertices incident from `iter`. That is,  $\min_j \{(iter, v_j) \in E\}$ . We prefer that your solution be a graph embedding either by hand-drawing and scanning or taking a picture, or drawing it digitally. If none of these options are available, you can simply list the set of edges in the graph.

**Answer:**



2. Develop an  $O(n^2)$  algorithm for solving the *longest path problem* in ordered, unweighted, directed graphs. Besides the functions defined above, you are free to use any function associated with graph data structure implementations.

**Answer:**

---

Algorithm 2: Longest path algorithm

---

```
def longest_path(v1):  
    length=0  
    for node in v1:  
        path=[node]  
        while node.hasOutgoingEdge():  
            node=node.getEdgeSmallestIndex()  
            path.append(node)  
        if len(path)>length:  
            length=len(path)  
  
    return length
```

---