

Data Structures and Algorithms Final Report 44251017

Huang Jiahui

Q1

(1) Describe an efficient algorithm that finds a maximum weight tour r . If you use an algorithm design paradigm, show its name.

The problem can be solved efficiently using Dynamic Programming on Trees. The core idea is to perform a single traversal of the tree (specifically, a post-order traversal or Depth First Search). For each node, we compute the longest path starting at that node and going downwards, while simultaneously tracking the longest path found anywhere in the tree so far.

Algorithm Design Paradigm: Dynamic Programming on a Tree.

```
1  // Assume a TreeNode structure exists with:
2  //   int weight;
3  //   List<TreeNode> children;
4
5  // --- Global variable ---
6  // Used to track the maximum weight found throughout the recursion.
7  int max_weight_so_far;
8
9
10 // --- Main Function ---
11 // Description: The entry point to start the algorithm.
12 FUNCTION findMaximumWeightTour(TreeNode root):
13     // Initialize the global maximum to a very small value.
14     max_weight_so_far = Integer.MIN_VALUE;
15
16     // Call the recursive helper to traverse the entire tree.
17     dfsHelper(root);
18
19     // Return the final recorded maximum weight.
20     RETURN max_weight_so_far;
21
22
23 // --- Recursive Helper Function (DFS) ---
24 // Description: Processes each node in a post-order fashion to calculate paths and
25 // update the global max.
26 // Returns: The weight of the longest "single-branch" path starting from 'node'
27 // and going downwards.
28 FUNCTION dfsHelper(TreeNode node):
29     // Base case: If the node doesn't exist, its contribution to any path is 0.
30     IF (node is null):
31         RETURN 0;
32
33     // Recursively call for all children.
```

```

32 List<Integer> child_down_paths = new ArrayList<>();
33 FOR EACH TreeNode child in node.children:
34     // Calculate the downward path weight from the child.
35     int path_from_child = dfsHelper(child);
36     // Add it to our list, but don't consider negative paths (treat them as
0).
37     child_down_paths.add( max(0, path_from_child) );
38
39     // Find the two longest downward paths coming from the children.
40     Sort child_down_paths in descending order;
41
42     int max1 = 0;
43     int max2 = 0;
44     IF (child_down_paths is not empty):
45         max1 = child_down_paths.get(0);
46     IF (child_down_paths has more than 1 element):
47         max2 = child_down_paths.get(1);
48
49     // Calculate the path "through" the current node and try to update the global
max.
50     int path_through_node = node.weight + max1 + max2;
51     max_weight_so_far = max(max_weight_so_far, path_through_node);
52
53     // Return the longest "single-branch" downward path from this node for its
parent to use.
54     RETURN node.weight + max1;

```

(2) Demonstrate how your algorithm finds maximum weight tour on the above tree.

We use the correct tree structure and perform a post-order traversal. Let `M` be `max_weight_so_far`, initialized to `Integer.MIN_VALUE`.

1. Leaves:

- `dfsHelper(40)` returns 40. `M = 40`.
- `dfsHelper(20)` returns 20. `M = 40`.
- `dfsHelper(5)` returns 5. `M = 40`.
- `dfsHelper(32)` returns 32. `M = 40`.
- `dfsHelper(22)` returns 22. `M = 40`.
- `dfsHelper(60)` returns 60. `M = 60`.
- `dfsHelper(3)` (leaf) returns 3. `M = 60`.

2. Internal Nodes (in post-order):

- Node 2: Child is 40. Returns `w(2) + 40 = 42`. `M = max(M, 42) = 60`.
- Node 10 (left): Children are 2, 20. `dfsHelper` returned 42 and 20.

- $\text{max1} = 42, \text{max2} = 20.$
- $\text{path_through_10(L)} = w(10) + 42 + 20 = 72. M = \max(60, 72) = 72.$
- Returns $w(10) + 42 = 52.$
- Node 10 (right): Child is 32. Returns $w(10) + 32 = 42. M = \max(M, 42) = 72.$
- Node 6: Children are 22, 60, 3(leaf). `dfsHelper` returned 22, 60, 3.
 - $\text{max1} = 60, \text{max2} = 22.$
 - $\text{path_through_6} = w(6) + 60 + 22 = 88. M = \max(72, 88) = 88.$
 - Returns $w(6) + 60 = 66.$
- Node 4: Children are 10(right), 6. `dfsHelper` returned 42, 66.
 - $\text{max1} = 66, \text{max2} = 42.$
 - $\text{path_through_4} = w(4) + 66 + 42 = 112. M = \max(88, 112) = 112.$
 - Returns $w(4) + 66 = 70.$
- Node 3 (root): Children are 10(left), 5, 4. `dfsHelper` returned 52, 5, 70.
 - $\text{max1} = 70$ (from node 4), $\text{max2} = 52$ (from node 10-left).
 - $\text{path_through_3} = w(3) + 70 + 52 = 125. M = \max(112, 125) = 125.$
 - Returns $w(3) + 70 = 73.$

The final value for `max_weight_so_far` is 125. This maximum was found when considering the root node 3 as the apex. The path is formed by combining the heaviest downward path from node 4 (4 -> 6 -> 60), the heaviest downward path from the left node 10 (10 -> 2 -> 40), and the apex node 3 itself.

The path is: 40 - 2 - 10(left) - 3 - 4 - 6 - 60. Weight = 40+2+10+3+4+6+60 = 125.

(3) Prove that your algorithm is correctly finding a tour having maximum weight.

Proof by Contradiction and Induction:

Let P be a maximum weight tour (path) in the tree T . Any path P has a unique "highest" node, which we call the *apex* of the path. The apex is the node on the path that is closest to the root of the tree.

Our algorithm iterates through every node u in the tree and considers it as a potential apex. For each u , it calculates the weight of the best possible path with u as its apex. The maximum of these values is the result.

Let u^* be the apex of the true maximum weight path P . The path P consists of u^* and at most two disjoint downward paths starting from children of u^* . Let these downward paths be P_1 (starting at child v_1) and P_2 (starting at child v_2). The total weight is $w(P) = w(u^*) + w(P_1) + w(P_2).$

The core of our algorithm is the function `dfsHelper(v)`, which we claim correctly computes the maximum weight of a path starting at v and descending into its subtree. We can prove this by induction on the height of the subtrees.

- Base case (height 0): For a leaf v , $\text{dfsHelper}(v)$ returns $w(v)$, which is correct.
- Inductive step: Assume $\text{dfsHelper}(c)$ is correct for all children c of v . Then the maximum downward path from v must consist of v itself plus the maximum downward path from one of its children c' . Our algorithm computes this as $w(v) + \max(0, \max_c(\text{dfsHelper}(c)))$. This is correct.

Since dfsHelper is correct, when our algorithm reaches the node u^* (the apex of the true max-weight path P), it will have correctly computed $\text{dfsHelper}(v_1)$ and $\text{dfsHelper}(v_2)$. These must be equal to $w(P_1)$ and $w(P_2)$. If they weren't, it would contradict that P is the overall maximum weight path.

Our algorithm calculates the candidate path weight at u^* as $w(u^*) + \max_1 + \max_2$, where \max_1 and \max_2 are the weights of the two heaviest downward paths from its children. This calculation will precisely equal the weight of P .

Since the algorithm performs this calculation for *every* node as a potential apex, it is guaranteed to perform it for the true apex u^* . Therefore, the final value of max_weight_so_far must be the weight of the maximum weight path P .

(4) Suppose that tree T has n nodes. Show the worst-case running time of your algorithm.

The running time of the algorithm is determined by the traversal and the work done at each node.

1. **Traversal:** The algorithm is based on a single Depth First Search (or post-order traversal) of the tree. In a tree with n nodes and $n-1$ edges, a DFS visits each node and each edge a constant number of times. This traversal has a complexity of $O(n + (n - 1)) = O(n)$.
2. **Work per Node:** For each node u , the algorithm performs the following operations:
 - It receives the return values from all its children.
 - It iterates through these values to find the top two largest. If node u has $\text{deg}(u)-1$ children (where $\text{deg}(u)$ is its degree), sorting or finding the top two takes time proportional to the number of children, which is $O(\text{deg}(u))$.
 - It performs a constant number of additions and comparisons.

3. **Total Time:** The total running time is the sum of the work done at each node during the traversal.

$$T(n) = \sum O(\text{deg}(u)) \quad , u \in V$$

In graph theory, the sum of the degrees of all vertices in a graph is equal to twice the number of edges:

$$\sum \text{deg}(u) = 2 | E |$$

For a tree with n nodes, the number of edges $|E|$ is $n-1$.

$$\sum \text{deg}(u) = 2(n - 1)$$

Therefore, the total time complexity is $O(2(n - 1)) = O(n)$.

The worst-case running time of the algorithm is $O(n)$, which is linear in the number of nodes.

Q2

Given recurrence:

$$T(n) = \begin{cases} c & n < 2 \\ T(n/2) + \log n & n \geq 2 \end{cases}$$

Expansion:

$$T(n) = T(n/2) + \log n$$

$$= T(n/4) + \log(n/2) + \log n$$

$$= T(n/8) + \log(n/4) + \log(n/2) + \log n$$

Continue expanding until $n/2^k < 2$, so $k = \log_2 n$.

Sum up all terms:

$$T(n) = c + \log n + \log(n/2) + \log(n/4) + \dots + \log 2$$

$$= c + \log n + (\log n - 1) + (\log n - 2) + \dots + 1$$

This is an arithmetic series, and the sum is:

$$T(n) = \Theta((\log n)^2)$$

Q3

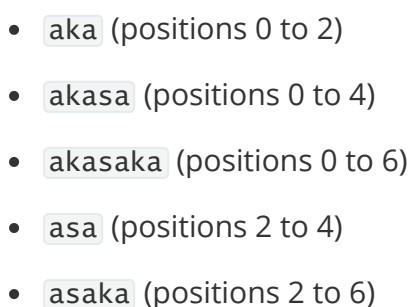
(1) Show the suffix trie for the following string: akasaka

To build a suffix trie, we first list all suffixes of the string. To ensure that each suffix has a unique path to a leaf, we append a special terminal character, \$, which is not in the alphabet. The string becomes akasaka\$.

The suffixes are:

1. akasaka\$ (index 0)
2. kasaka\$ (index 1)
3. asaka\$ (index 2)
4. saka\$ (index 3)
5. aka\$ (index 4)
6. ka\$ (index 5)
7. a\$ (index 6)

By inserting these suffixes into a trie, we get the following structure. The numbers in the leaf nodes represent the starting index of the corresponding suffix.



- aka (positions 4 to 6)

(3) Show an efficient algorithm to answer query aXa , using a suffix trie. Here, the answers can be the start and end positions of the given string.

This algorithm solves the general query $c1 \times c2$, which finds all substrings starting with character $c1$ and ending with character $c2$. The query aXa is a special case where $c1$ and $c2$ are the same.

The algorithm leverages a suffix trie to efficiently perform the first part of the task: identifying all possible starting positions.

Algorithm Strategy

The algorithm is divided into two main phases:

1. **Find Start Positions:** The suffix trie groups all suffixes by their prefixes. This means all suffixes of the string that start with $c1$ will share a common path from the root. We can exploit this to find all occurrences of $c1$ in a single traversal of the relevant subtree.
2. **Find End Positions:** For each starting position i found in the first phase, we perform a linear scan through the rest of the string to find all positions j (where $j > i$) that contain the character $c2$.

```
1 // Assume the Suffix Trie for the string S is pre-built.
2
3 FUNCTION Find_Pattern_Matches(S, trie, c1, c2):
4     // Inputs:
5     //   S: The original string of length n.
6     //   trie: The pre-built suffix trie for S.
7     //   c1: The starting character of the pattern.
8     //   c2: The ending character of the pattern.
9     // Output:
10    //   A list of (start, end) position pairs for all matching substrings.
11
12    // --- Phase 1: Find Start Positions using the Suffix Trie ---
13
14    // 1. Traverse from the root to the node corresponding to the first character.
15    Node c1_node = trie.root.getChild(c1)
16
17    // 2. If no such path exists, there are no matches.
18    IF c1_node is NULL:
19        RETURN empty list
20
21    // 3. Collect all start indices from the leaves of the subtree rooted at
22    //    c1_node.
23    //    This gives us the 0-indexed positions of all occurrences of c1 in S.
24    List<Integer> start_positions = trie.collectAllLeafIndicesFrom(c1_node)
25
26    // --- Phase 2: Find End Positions by Scanning ---
27
28    // 4. Initialize a list to store the results.
29    List<Pair<Integer, Integer>> results = new empty list
```

```

29
30 // 5. For each valid starting position...
31 FOR EACH start_pos IN start_positions:
32     // 6. ...scan the rest of the string to find valid ending positions.
33     FOR j FROM (start_pos + 1) TO (n - 1):
34         IF S[j] == c2:
35             // 7. A match is found. Add the (start, end) pair to the results.
36             results.add(Pair(start_pos, j))
37
38 // 8. Return the complete list of matching positions.
39 RETURN results

```

(4) Suppose the input string has length n , and d distinct characters are used. What is the time complexity of your algorithm in (3)?

The time complexity of the algorithm is determined by two main operations during the query phase:

1. Finding Start Positions: This involves traversing a subtree within the suffix trie. In the worst case for a standard (uncompressed) trie, the size of this subtree can be proportional to n^2 , leading to a complexity of $O(n^2)$ for this step.
2. Finding End Positions: For each of the k start positions found, the algorithm scans the remainder of the string. This nested loop structure results in a complexity of $O(k \cdot n)$.

In the worst-case scenario, k can be as large as n (e.g., for the string "aaaaa..."). Therefore, the total time complexity is dominated by these operations, resulting in $O(n^2 + n \cdot n) = O(n^2)$. The number of distinct characters, d , does not affect the asymptotic time complexity.