COL351 Assignment 1

Mallika Prabhakar, 2019CS50440 Sayam Sethi, 2019CS10399

August 2021

Contents

1	Question 1 (MSTs)1.a Unique MST1.b Algorithm Sketch	
2	Question 2 (Huffman Encoding)2.a Optimal Huffman Encoding	6 6 9
3	Question 3 (Greedy Approaches)3.a Optimal Selection3.b Optimal Seating	

1 Question 1 (MSTs)

Let G be an edge-weighted graph with n vertices and m edges satisfying the condition that all the edge weights in G are distinct.

1.a Unique MST

Question 1.a

Question. Prove that G has a unique MST.

Proof. We will prove this by induction on the size of G using an idea similar to Kruskal's algorithm discussed in the class.

Hypothesis:

$$h(n): \forall G = (V, E): |V| = n \implies MST(G) \text{ is unique}$$
 (1)

Base case: n=1 is true since there is no edge and $MST(G)=(V,\phi)$ is unique.

```
Induction Step: Assume h(n-1) is true for n \ge 2, now for h(n): (Note: This proof assumes each edge to be an unordered pair of vertices)
```

Consider Kruskal's algorithm,

Algorithm 1 Recursive MST Routine — Kruskal's algorithm

```
1: procedure MST(G)
 2:
        e_0 \leftarrow (x, y) be edge with least weight
        H \leftarrow G
 3:
        remove x, y from H and add new vertex z
 4:
        for all v such that v is neighbour of x or y do
 5:
            add (v,z) to H
 6:
            wt(v,z) \leftarrow \min(wt(v,x), wt(v,y))
 7:
           if wt(v,x) < wt(v,y) then
 8:
               map(v,z) \leftarrow (v,x)
 9:
            else
10:
                map(v,z) \leftarrow (v,y)
11:
            end if
12:
        end for
13:
        T_H \leftarrow MST(H)
14:
        T_G \leftarrow (V, \{e_0\})
15:
        for all e \in T_H do
16:
            if e is not incident on z then
17:
                add e to T_G
18:
            else
19:
                add map(e) to T_G
20:
21:
            end if
        end for
22:
        return T_G
23:
24: end procedure
```

In the above algorithm, it is clear that H has n-1 vertices. Thus, by our assumption, h(n-1) is true and hence T_H is unique. Also, we know that T_G is a valid MST, from the correctness of Kruskal's algorithm. Now, assume by contradiction that T_G is not unique. Then there exists an MST, say $T' \neq T_G$.

Claim 1.1. e_0 cannot be in T'

Proof. This is because, if e_0 were in T', then $T \setminus \{e_0\} \neq T' \setminus \{e_0\}$ and thus, there would be two different MSTs for H which would be a contradiction to our assumption. Thus, $e_0 \notin T'$.

Consider the path from x to y in T'. Since $e_0 = (x, y)$ is not present in T', there exists a different path, say $P = (f_1, f_2 \cdots, f_k)$ where $f_i \in E(T'), 1 \le i \le k$. We know that

 $wt(f_i) > wt(e_0), 1 \le i \le k.$

Swap any of the f_i with e_0 and let the subgraph formed be T'', i.e., $T'' = T' \setminus \{f_i\} \cup \{e_0\}$. We know T'' is a spanning tree of G since V(T'') = V(G) and there are no cycles formed on performing the swap operation (this can be proven using contradiction as discussed in the lecture).

Now, consider the weight of T'':

$$wt(T'') = wt(T') - wt(f_i) + wt(e_0)$$

$$\implies wt(T'') < wt(T')$$
(2)

We have shown that the total weight of T'' is lesser than the weight of T'. However, this is a contradiction to the fact that T' is the MST of G. Thus our assumption that T_G is not the unique MST of G was wrong. Therefore, h(n) is true.

This completes the induction and the proof that if all edge weights in a graph are distinct, then its MST is unique.

1.b Algorithm Sketch

Question 1.b

Question. If it is given that G has at most n+8 edges, then design an algorithm that returns a MST of G in O(n) running time.

Solution. The idea is to use the previous result along with the fact that the number of edges to be removed to form a spanning tree is m - (n - 1) which is at most (n + 8) - (n - 1) = 9, assuming that G was initially connected (else no MST exists). The algorithm is as follows:

Algorithm 2 Compute MST for 1.b

```
1: procedure MST(G)

2: while |E(G)| > |V(G)| - 1 do

3: C \leftarrow findCycle(G)

4: e \leftarrow edge with largest weight in C

5: remove e from G

6: end while

7: return G

8: end procedure
```

The procedure findCycle calls a DFS function on G which uses graph colouring and returns the first cycle it finds:

Algorithm 3 findCycle

```
1: procedure FINDCYCLE(G)
         v \leftarrow \text{any vertex of } G
 3:
         colour \leftarrow map of vertices initialised to zero
         parent \leftarrow map of vertices initialised to null
 4:
         (u, v) \leftarrow \mathrm{DFS}(G, v, \mathrm{colour}, \mathrm{parent}, \mathrm{null})
 5:
                           ▷ DFS returns the bottommost and topmost vertex of the cycle
 6:
         if DFS returned null then
 7:
             return null
 8:
         end if
 9:
         C \leftarrow \text{empty array of edges}
10:
         add (u, v) to C
11:
         while u \neq v do
12:
             add (u, parent(u)) to C
13:
             u \leftarrow \operatorname{parent}(u)
14:
         end while
15:
16:
         return C
17: end procedure
```

The *DFS* function looks as follows:

Algorithm 4 Identify cycle using colouring and DFS

```
1: procedure DFS(G, v, \text{colour}, \text{parent}, p)
         parent(v) \leftarrow p
 2:
         \operatorname{colour}(v) \leftarrow 1
 3:
 4:
         for all u such that u is neighbour of v in G do
             if colour(u) is 2 then
 5:
                                                              \triangleright there is a forward edge from v to u
                  return (u, v)
 6:
             else if colour(u) is 0 then
                                                                                           \triangleright u is unvisited
 7:
                  value \leftarrow dfs(G, u, colour, parent, v)
 8:
                  if value is not null then
                                                                         \triangleright cycle found in subtree of u
 9:
                      return value
10:
                  end if
11:
             end if
12:
         end for
13:
         \operatorname{colour}(v) \leftarrow 2
14:
         return null
                                                                     \triangleright no cycle found in subtree of v
15:
16: end procedure
```

We will assume without proof that the DFS (Algorithm 4) function is correct and it takes O(n+m) time since the algorithm is a standard one and has been discussed in the lecture.

Now consider the function findCycle (Algorithm 3), lines 3, 4 take O(n) time and line

5 takes O(n+m) time. The while loop (lines 9–12) traverses up from u to v and each iteration takes O(1) time. Therefore the entire while loop completes in O(n) time (since the graph has n vertices and hence the loop cannot run for more than n iterations). Therefore the total time complexity of findCycle is O(n+m).

We will now prove termination and compute complexity of Algorithm 2, which contains the code for computing MST:

Termination: The while loop terminates when |E| = |V| - 1, that is, when the graph is a tree (since it assumes that the graph is connected). In each iteration of the while loop, findCycle returns a valid cycle since |E| > |V| - 1 and the graph is connected. After having found the cycle, we remove the edge with largest weight from G and therefore |E| reduces by 1. Since |V| remains constant, the while loop terminates after a finite number of steps.

Time Complexity: The while loop runs for exactly m - (n - 1) iterations, which is O(m - (n - 1)) = O((n + 8) - (n - 1)) = O(1) for the given constraints. Each iteration of the while loop calls findCycle which runs in O(n + m) = O(n + (n + 8)) = O(n). Finding the edge with least weight is O(n) since a cycle cannot have more than n edges. Removing this edge from G can be implemented in as worse as O(n) (better implementations in O(1) and $O(\log(n))$ time exist but this won't change the complexity of the algorithm as we will show). Thus, each iteration of the while loop takes O(n) time and the total time complexity of Algorithm 2 is:

$$T(MST) = O(\text{iterations of while loop} \times \text{complexity of each iteration})$$

= $O(O(1) \times O(n)) = O(n)$ (3)

Correctness: We now proceed to prove the correctness of the algorithm, using the following claim,

Claim 1.2. If a cycle has edges of distinct weights, the edge with the largest weight can not be a part of any MST

Proof. Let us assume by contradiction that the claim is false, then there exists an MST, say T such that the largest edge of cycle C (with distinct weights) is present in T. Let that edge be e = (x, y). Now consider the paths from x to y in G. There exists at least another path from x to y, which is exactly equal to $C \setminus \{e\}$, call it P. Consider the edge in P which is not in T, say f = (p, q). We know such an edge exists since T is acyclic. Now, consider $T' = T \setminus \{e\} \cup \{f\}$. We will now prove that T' is a spanning tree using the following claim:

Claim 1.3. Consider any edge m = (a, b) in G which is not in T (spanning tree of G). Let n = (u, v) be any edge on the unique path from a to b in T. Then on swapping m with n in T, we get another spanning tree of G.

Proof. If $path_T(u, a)$ does not exist in T, then swap u, v (for ease of notation). Consider the graph $T \setminus \{n\}$. Define set

$$S = \{(c,d) \mid path_T(c,d) = \{k_1, k_2, \dots, u, v, \dots, k_l\}\}$$
(4)

All pair of vertices in this set are disconnected since all paths in the tree are unique. Now, consider the path

$$P_{T'} = \{c = k_1, k_2, \dots, u\} \cup path_T(u, a) \cup path_T(b, v) \cup \{\dots, d = k_l\}, \forall (c, d) \in S \quad (5)$$

Now, consider the graph $T' = T \setminus \{n\} \cup \{m\}$. All paths given by $P_{T'}$ are present in T' and thus, all pairs of vertices in S are connected in T'. Since all other paths are the same in T and T', T' is connected. Since |E(T')| = |V(T')| - 1, T' is a tree and also a spanning tree of G. This completes the proof of the claim.

Thus, from Claim 1.3, we know that T' is an MST of G. Consider the weight of T':

$$wt(T') = wt(T) - wt(e) + wt(f)$$

$$\implies wt(T') < wt(T), \text{ since } wt(e) > wt(f)$$
(6)

This is a contradiction to the fact that T is an MST of G. Therefore our assumption that Claim 1.2 is incorrect was wrong. This proves the correctness of Claim 1.2.

Now consider Algorithm 2. In each iteration of the algorithm, we remove the largest edge of a cycle from G. Let the new graph obtained be G'. Therefore our algorithm transforms the problem from G to G'. We need to show that both graphs have the same MST.

From Claim 1.2, we know that e cannot be in any MST of G and from Question 1.a, we know that the MST of G will be unique. Therefore, the MST of G and $G' = G \setminus \{e\}$ will be the same. This completes the proof of correctness of Algorithm 2.

2 Question 2 (Huffman Encoding)

2.a Optimal Huffman Encoding

Question 2.a

Question. What is the optimal binary Huffman encoding for n letters whose frequencies are the first n Fibonacci numbers? What will be the encoding of the two letters with frequency 1, in the optimal binary Huffman encoding?

Solution. We begin by observing the property of Fibonacci numbers:

$$f_n = f_{n-1} + f_{n-2} \ \forall n \ge 3$$

and, $f_1 = f_2 = 1$ (7)

We are given an alphabet $A = (a_1, a_2, ..., a_n)$ such that it has a frequency vector $F = (f_1, f_2, ..., f_n)$. Before finding the encoding, consider the sum of first k Fibonacci

numbers, call it s_k :

$$s_{k} = f_{1} + f_{2} + \dots + f_{n-2} + f_{n-1} + f_{n}$$

$$\implies s_{k} = f_{1} + f_{2} + \dots + f_{n-2} + f_{n+1}$$

$$\implies s_{k} = s_{k-2} + f_{n+1}$$

$$\implies s_{k} - s_{k-2} = f_{n+1}$$
(8)

On performing telescopic summation over Equation 8 (for k > 2), we get the following:

This Equation 9 takes a form similar to Equation 7 and thus, $s_k + 1 = f_m$ for some m. On substituting value of k = 1:

$$s_{1} + 1 = f_{m}$$

$$\Rightarrow f_{m} = 2$$

$$\Rightarrow m = 3$$

$$\Rightarrow s_{k} + 1 = f_{k+2}$$

$$\Rightarrow s_{k} = f_{k+2} - 1$$
(10)

Now consider the Huffman tree for |A| = n. Each of the frequency f_i $(1 \le i \le n-2)$ is less than f_n and sum of all frequencies f_i $(1 \le i \le n-2)$, i.e., $s_{k-2} = f_n - 1$ is less than f_n . We also know that a_i is merged at the same time or before a_j for any i < j. From this, we can formulate the merging strategy with the help of the following inductive claim:

Claim 2.1. The optimal Huffman tree for A with frequency vector F is constructed in a way such that (a_1, a_2, \ldots, a_i) is merged in the first i-1 steps $\forall i: 1 \leq i \leq n$.

Proof. Base case: i = 1 is trivially true since a_1 is a leaf node and is *merged* in 0 merges.

Induction Step: Assume the claim is true for i-1. After i-2 merges, $(a_1, a_2, \ldots, a_{i-1})$ have been merged into $tree(a_1, a_2, \ldots, a_{i-1})$, and the frequency vector will be as follows,

$$F = (f_1 + f_2 + \dots + f_{i-1}, f_i, f_{i+1}, \dots, f_n)$$

$$F = (s_{i-1}, f_i, f_{i+1}, \dots, f_n)$$

$$F = (f_{i+1} - 1, f_i, f_{i+1}, \dots, f_n), \text{ from Equation } 10$$
(11)

It is easy to see that the least two frequencies in the frequency vector are f_i , $f_{i+1} - 1$ which correspond to a_i and $tree(a_1, a_2, \ldots, a_{i-1})$. Therefore the $(i-1)^{th}$ merge will merge these two into $tree(a_1, a_2, \ldots, a_i)$.

We have shown that a_i is merged in the $(i-1)^{\text{th}}$ step and from induction we know that $(a_1, a_2, \ldots, a_{i-1})$ are merged before (i-1) steps and thus, (a_1, a_2, \ldots, a_i) are merged in i-1 steps. This completes the induction and proves the claim.

Therefore, from Claim 2.1, we know that a_n is merged in the last step (which is the $(n-1)^{\text{th}}$ step) and hence it is encoded using a single bit. We can now inductively define the encoding for each alphabet (for n > 1):

Claim 2.2.
$$a_i$$
 is encoded as $\underbrace{11...1}_{n-i \text{ times}} 0$ for $n \ge i > 1$ and a_1 is encoded as $\underbrace{11...1}_{n-1 \text{ times}}$

Proof. For i > 1, we will prove the claim using induction.

Base case: From Claim 2.1, we know that a_n will be merged in the last step and thus it is encoded using a single bit, we can choose this bit to be 0 and thus $enc(a_n) = \underbrace{11...1}_{0} 0 = 0$ and the claim is true for n.

Induction Step: Assume the claim is true for i+1, i.e., $enc(a_{i+1}) = \underbrace{11...1}_{n-(i+1) \text{ times}} 0$.

From the proof of the previous claim, we know that a_{i+1} and $tree(a_1, a_2, ..., a_i)$ are siblings and thus, the encoding of the root of $tree(a_1, a_2, ..., a_i)$ will be $\underbrace{11...1}$.

From the base case, we know that a_n is encoded using a single bit with respect to the root of the tree. Therefore, with respect to $tree(a_1, a_2, \ldots, a_i)$, we know that a_i is encoded using a single bit. Let that bit be 0. We then have the complete encoding of a_i as:

$$enc(a_i) = enc(tree(a_1, a_2, \dots, a_i)).0$$
 (. denotes concatenation)
= $\underbrace{11 \dots 1}_{n-i \text{ times}} 0$ (12)

This completes the induction for i > 1 and we now show the correctness of the claim for i = 1.

We know that a_1 and a_2 are siblings in the Huffman tree and thus they differ in their representation in exactly the last bit. Therefore, $enc(a_1) = \underbrace{11 \dots 1}_{n-1 \text{ times}}$. This completes

the proof of the claim.
$$\Box$$

Thus, we have computed the optimal Huffman encoding for the alphabet $A = (a_1, a_2, \ldots, a_n)$ which has frequency vector as $F = (f_1, f_2, \ldots, f_n)$ and we restate Claim 2.2:

In the optimal Huffman encoding for A with frequency F such that |A| = n, a_i is encoded as $\underbrace{11\ldots 1}_{n-i \ times} 0$ for $n \geq i > 1$ and a_1 is encoded as $\underbrace{11\ldots 1}_{n-1 \ times}$ (and for n = 1,

$$a_n = a_1 = 0$$
 trivially).

2.b Prove no Efficiency

Question 2.b

Question. Suppose you aim to compress a file with 16-bit characters such that the maximum character frequency is strictly less than twice the minimum character frequency. Prove that the compression obtained by Huffman encoding, in this case, is same as that of the ordinary fixed-length encoding.

Proof. Instead of proving the solution for a 16 - bit encoding, we will prove this for any k - bit encoding with the help of the following claim:

Claim 2.3. For an encoding of an alphabet A_k of size 2^k such that $f_1 < 2f_{2^k}$ where F is the frequency vector in increasing order, a_{2m-1} and a_{2m} are siblings for $1 \le m \le 2^{k-1}$ and the same inequality holds for the parents of a_i $(1 \le i \le 2^k)$

Proof. Let $2^k = n$. Consider the frequency vector F, the least two frequencies are of a_1 and a_2 . Let them be merged into b_1 where frequency of b_1 is $g_1 = f_1 + f_2$. It is easy to see that g_1 is greater than all f_i $(1 \le i \le n)$ using the fact that $f_n < 2f_1 < g_1$. We will now prove the following hypothesis using induction:

$$h(i): b_p = parent(a_{2p-1}, a_{2p}), 1 \le p \le i, \text{ and } g_p < g_q, 1 \le p < q \le i$$
 (13)

Base case: We have shown the base case, h(1) to be true above.

Induction Step: Assume it is true for i-1, and consider the frequency vector after these i-1 merges in sorted order,

$$F = (f_{2i-1}, f_{2i}, \dots, f_{n-1}, f_n, g_1, g_2, \dots, g_{i-1})$$
(14)

Thus, the least two frequencies are of a_{2i-1} and a_{2i} and these are merged into $b_i = parent(a_{2i-1}, a_{2i})$. Now,

$$g_i = f_{2i-1} + f_{2i}$$

we know that,

$$f_{2i-1} > f_{2(i-1)-1}$$
, and $f_{2i} > f_{2(i-1)}$
 $\implies f_{2i-1} + f_{2i} > f_{2(i-1)-1} + f_{2(i-1)}$
 $\implies g_i > g_{i-1}$ (15)

Therefore, h(i) is also true and this completes the induction Now, consider the statement of h(n/2) (n/2 is an integer for k > 0):

$$b_p = (a_{2p-1}, a_{2p}), 1 \le p \le n/2 = 2^{k-1}$$
(16)

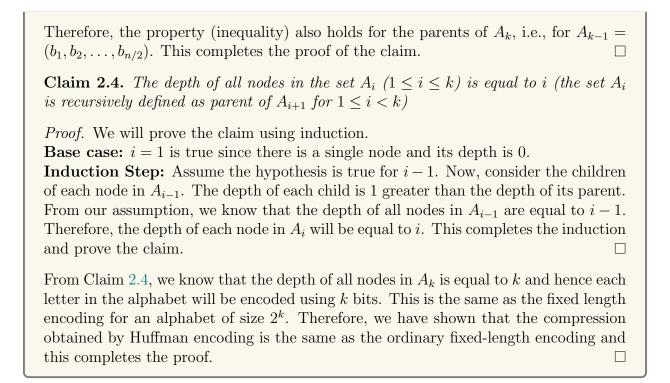
This completes the proof of the first half of the claim, now consider g_1 and $g_{n/2}$,

$$2a_1 > a_n$$
, from the property of the alphabet

$$\implies 4a_1 > 2a_n$$

$$\implies 2(a_1 + a_2) > a_n + a_{n-1}$$

$$\implies 2g_1 > g_{n/2}$$
(17)



3 Question 3 (Greedy Approaches)

3.a Optimal Selection

Question 3.a

Question. Present an efficient algorithm which outputs best choice of party invitees as per following specifications:

- 1. Input: list of n people and list of pairs who know each other (undirectional)
- 2. Every person invited should have at least k people they know and k people they don't know at the party
- 3. k is set to be 5 in the question but we used k for generality

Solution. The main idea here is to obtain an adjacency list adj and people which contains information about vertices and their degrees. We keep a variable done which tells us if the people set is being changed or hasn't changed after a complete iteration of while loop. If there is no change i.e. done = true, the while loop breaks. It can also break when people becomes empty. Otherwise it simply denotes that guests are still being removed to find the optimal invitee list.

Algorithm 5 Generate list of invitees for 3.a

```
1: procedure Invite(G)
 2:
        initialise adj
                         ⊳ empty hashmap of balanced BSTs created for adjacency list
 3:
        initialise people
                                            ▶ empty balanced BST of degree, vertex pairs
        for all edge e = \{u, v\} in E(G) do
                                                                                    ⊳ sets adi
 4:
           add v to adj[u]
 5:
           add u to adj[v]
 6:
        end for
 7:
        for all vertex v in V(G) do
 8:
                                                                                ⊳ sets people
           add \{degree(v), v\} to people
 9:
        end for
10:
        done \leftarrow false
                                                            \triangleright denotes if people is changing
11:
        while not done do
12:
           done \leftarrow true
13:
           if people is empty then
14:
               break
15:
16:
           end if
17:
           \{d,v\} \leftarrow \text{vertex with minimum degree in } people
           if d < k then
18:
               done \leftarrow false
19:
               deletePerson(v)
20:
           end if
21:
           if people is empty then
22:
               break
23:
           end if
24:
           \{d,v\} \leftarrow \text{vertex with maximum degree in } people
25:
           if size(people) - d - 1 < k then
26:
               done \leftarrow false
27:
               deletePerson(v)
28:
29:
           end if
30:
        end while
        initialise invitees
31:
                                                                               ⊳ empty array
        for all {degree, person} in people do
32:
33:
           add person to invitees
        end for
34:
        return invitees
35:
36: end procedure
```

The procedure deletePerson takes the vertex u into consideration for removal as a parameter and updates the adjacency list adj and people. Idea behind the algorithm is to look at all the neighbours of v and remove the edge $\{u, v\}$ from adj and update degree of v in people. If degree of v becomes zero, don't add it again.

Since people is a balanced BST, time complexity of removal and insertion is O(log|V|)

because size(people) = |V|. Also, adjacency list can be implemented as a hashmap of balanced BSTs, so locationg a vertex v takes O(1) time and removing a vertex u from the hashmap[v] takes $O(log\ degree)$

Total time complexity of the procedure deletePerson hence becomes

```
O(degree(u)) \times O(\log|V| + \log(degree(v))) + O(\log|V|) = O(degree(u) \times \log|V|) 
(18)
```

```
Algorithm 6 sub-algorithm for deletePerson
```

```
1: procedure DELETEPERSON(u)
 2:
       for all vertex v in neighbours of u do
          remove \{degree(v), v\} from people
 3:
          remove u from adj(v)
 4:
          if degree(v) \neq 0 then
 5:
              insert \{degree(v), v\} into people
 6:
           end if
 7:
       end for
 8:
 9:
       remove \{degree(u), u\} from people
10:
       return
11: end procedure
```

Initialisation

adj and people are initialised as an empty hashmap and set. We then use a for loop over E(G) to set values for adj. Then adj is used to set values for people using a for loop over V(G) as {degree, vertex} pairs. done is set to false initially and is used as a measure to check if further changes are being made to people or not.

While loop works on the condition of while *done* is not true. Inside the loop, *done* is set to true and If *people* turns out to be empty, it breaks. Different conditions are checked which either modify *people* and set *done* to false or make no change.

Maintenance

For every iteration of while loop, done is set to true to observe which can be changed back to false if a new change occurs. After checking that people is non empty, we look at the pair with minimum degree $\{d,v\}$. If the d < k, it implies v knows less than k people which is unfavourable. Hence we perform deletePerson(*) operation on v. Again, we check if people is now empty. If it is, then we break out of the while loop, else continue. Now we look at the pair with maximum degree $\{d,v\}$. If size(people) - d - 1 < k it implies v doesn't know less than k people which is again, unfavourable. Hence we perform deletePerson(*) operation on v and move on to next iteration.

Termination

The algorithm terminates if the while loop terminates. The while loop terminates if either *people* becomes empty or done stays true i.e. no new deletePerson(*) operation was performed. The size of *people* keeps decreasing until one of the two conditions mentioned above happens.

Case 1: people is empty

There must be a *deletePerson()* operation performed in the last iteration else the algorithm would have terminated then. Now, the size of *people* is 0. Since there are no people to invite to the part anymore, while loop breaks based on first if condition. It terminates.

Case 2: done is false

There was a *deletePerson()* operation performed in the last iteration for reason same as case 1. *done* is set to true. Now, if the minimum or maximum degree vertex was not in accordance to question statement, it will be removed and *done* is set to false again and this will repeat until Case 1 or Case 3 is obtained since size(people) will keep on decreasing after every iteration according to our logic.

Case 3: done is true

After the deletePerson() operation in last iteration, in current iteration, done is set to true. It is observed that both minimum and maximum degree follow the condition specified in the question. Since no if statement is valid, done is not set to false. Before the next iteration, the while loop condition is check and since done was true, not done is false and while loop terminates.

After termination of loop, a list invitees is initialised as an empty list. for every pair in people, the person v is added to invitees if people is not empty.

invitees is returned

Proof of correctness

To prove the correctness of our algorithm, we will show the following:

- 1. Removing nodes with degree $\langle k \rangle$ is necessary
- 2. Removing nodes with degree > p k 1 is necessary (where p is the current number of vertices in the graph)
- 3. No other removal is needed

The above three are necessary and sufficient to prove the correctness of our solution.

1. Proof of 1

Consider a vertex v that has a degree

Consider a vertex v that has a degree < k at any point of time. Let its degree be d. For any future configuration (selection) of vertices from the graph that

contains v, the degree of v cannot be larger than d. This is because no new edges are added to the graph and edges are only removed, which might include neighbours of v.

Therefore, it is necessary to remove nodes with degree < k.

2. Proof of 2

Consider a vertex v that has a degree > p - k - 1 at any point of time. Let its degree be d. For any future configuration (selection) of vertices from the graph that contains v, the degree of v will be given by d - r, where r is the number of neighbours of v removed. Now, number of vertices in the graph that have been removed upto this configuration is t(> r). Therefore, the number of people v doesn't know is equal to,

$$(p-t) - (d-r) - 1 = (p-d+1) - (t-r) < p-d+1 < k$$
(19)

Therefore, the number of people v doesn't know can never increase in any future configuration and hence removing v is necessary.

3. Proof of 3

We have proved that the removals that we have done in 1 and 2 are necessary. We will not prove the sufficiency of the same. Consider any vertex that is present in G after all necessary removals. Each vertex v in G knows at least k other people in G else 1 would not have completed. Each vertex v also does not know at least k people since 2 has completed. Therefore every vertex v in G satisfies the given constraints in the problem and hence the necessary removals are also sufficient.

This completes the proof of correctness of the algorithm.

Time complexity

Time complexity of deletePerson(u) is $O(degree(u) \times \log |V|)$

For the algorithm,

Setting adj takes O(|E|) time

Setting people takes O(|V|) time

Setting invitees takes O(|V|) time

finding minimum and maximum degree vertices takes $O(\log |V|)$ time

In the worst case, while loop iterates over all the vertices (O(|V|)) iterations)

Total time complexity of the algorithm can be observed as follows:

$$O(|E|) + O(|V|) + O(\text{total operations inside while loop}) + O(|V|)$$

$$= O(|E|) + O(|V|) + O(|V| \log |V| + \sum_{u \in V} O(degree(u)) \times \log |V|) + O(|V|)$$

$$= O(|E| \times \log |V|)$$

$$= O(n^2 \log n)$$
(20)

Where n is |V| (given) and |E| is n^2 in worst case scenario

3.b Optimal Seating

Question 3.b

Question. Suppose finally Alice invited n_0 out of her n friends to the party. Her next task is to set a minimum number of dinner tables for her friends under the constraint that each table has a capacity of ten people and the age difference between members of each dining group should be at most ten years. Present a greedy algorithm to solve this problem in $O(n_0)$ time assuming the age of each person is an integer in the range [10,99].

Solution. We have to implement an $O(n_0)$ greedy approach for the fore-mentioned question. Idea behind our algorithm is to create a frequency array which stores how many people of a certain age are there in the invited group from age 10 to 100. We then start grouping 10 people together from the smallest to largest age. If the difference between youngest and oldest person exceeds 10, we stop the group and create a new table. Following is the algorithm for the same:

Algorithm 7 Find number of tables needed

```
▶ Finds the minimum number of tables needed
 1: procedure Seat(people)
        arr[100] \leftarrow initialised as zero array
                                                                    ⊳ index 0 to 9 aren't used
 3:
        for person in people do
                                                                     ⊳ get the frequency array
            arr[age(person)] \leftarrow arr[age(person)] + 1
 4:
        end for
 5:
        set index, smallest, tcount \leftarrow 0
                                                             ▶ helpful intermediate variables
 6:
        set toAdd \leftarrow 10
 7:
        while index < 100 \text{ do}
 8:
            if smallest = 0 then
                                            ▷ smallest=0 signifies no person on a table yet
 9:
                smallest \leftarrow index
10:
11:
            end if
            if arr[index] = 0 then
                                                                       ⊳ no person of that age
12:
                index \leftarrow index + 1
13:
            else if index > smallest + 10 then
14:
                                                                                 ⊳ large age gap
                tcount \leftarrow tcount + 1
15:
16:
                smallest \leftarrow 0
                toAdd \leftarrow 10
17:
            else if arr[index] > toAdd then \Rightarrow table occupancy can reach atmost 10
18:
                arr[index] \leftarrow arr[index] - toAdd
19:
                tcount \leftarrow tcount + 1
20:
                toAdd \leftarrow 10
21:
                smallest \leftarrow 0
22:
            else
23:
                toAdd \leftarrow toAdd - arr[index]
24:
                arr[index] \leftarrow 0
25:
                index \leftarrow index + 1
26:
            end if
27:
        end while
28:
29:
        if toAdd \neq 10 then
                                                     ▶ take last non filled table into account
30:
            tcount \leftarrow tcount + 1
        end if
31:
        return tcount
32:

    ▶ total number of tables required is returned

33: end procedure
```

Now we shall prove the correctness of the algorithm:

Initialisation

We generate a frequency array by declaring a zero array of size 100 and incrementing the value stored at an index by 1. The update at an index is done based on what age value the for loop comes across. Updated index is same as age.

index variable is used to refer the current index in the while loop. toAdd variable

shows how many more people are to be added at the recent table. smallest denotes the smallest value of age of people sitting at a table. smallest = 0 means that the table is empty. tcount denotes the table number. All of the forementioned variables are initialised to zero except toAdd which is initialised at 10.

A while loop is run. The termination condition depends on the index variable. The loop keeps incrementing index until it finds the index at which arr[index] is not 0.

Maintenance

For maintenance, we look at the iterations of the while loop.

smallest = 0:

smallest = 0 implies that the table is new and there are no people seated on it.

The following cases arise in the if-else block:

Case 1: arr[index] is 0

Here, we move to the next index since there are no people of this age invited to the party.

Case 2: age gap between youngest person at table and index is > 10

Here, we cannot add the people at this index to the table since it violates our condition of the maximum age gap at table to be 10. We increase tcount by 1 and reset values of smallest and toAdd to 0 and 10 respectively.

Case 3: number of people with age same as index are more than toAdd

In this case, the current table can be filled completely. We subtract the number of people to be added from arr[index] and increment tcount by 1. toAdd and smallest are reinitialised.

Case 4: number of people with age same as index are $\leq toAdd$

Number of people of age same as index is not enough to fill the current table. We still allot them the table and move to next index. toAdd is decremented by arr[index]. arr[index] is set to 0.

Since last group of statements is an else and no other conditions remain, our maintenance step is complete.

Termination and Time Complexity

In each iteration of the while loop, the value of index increases by 1 or the total sum of arr decreases by some non-zero number or tcount increments by 1. Thus, the maximum number of iterations of the loop is bounded by:

$$100 + sum(arr) + \max(tables) = 100 + O(n_0) + O(n_0)$$

= $O(n_0)$ (21)

Therefore, the while loop runs for $O(n_0)$ iterations and each iteration takes O(1) time since it involved making finite checks and updates all of which take O(1) time.

The loop from lines 3-5 takes $O(n_0)$ time since it involves iterating over the entire

array of size n_0 once.

Therefore, the total time complexity of the algorithm is $O(n_0)$.

Proof of Correctness

To prove the correctness of the algorithm, we first prove the following claim:

Claim 3.1. If there exists an optimal seating where p_i, p_k are on one table and p_j is on another table such that $age(p_i) < age(p_j) < age(p_k)$, there exists another optimal seating where p_i, p_j are on one table and p_k is on another table.

Proof. Since p_i and p_k are on the same table, we know that $age(p_k) \leq age(p_{l1}) + 10$, where p_{l1} is the person with least age on the table. Consider the person with smallest age on the table containing p_j , say p_{l2} . We know that $age(p_j) \leq age(p_{l2}) + 10$. Therefore, we have that:

$$age(p_k) \le age(p_{l2}) + 10$$
, and $age(p_i) \le age(p_{l1}) + 10$ (22)

Therefore, it is possible to swap the positions of p_k and p_j . Thus, we have obtained an arrangement where p_i, p_j are on the same table and p_k is on a different table. This completes the proof of the claim.

Consider the person with the smallest age, say p_0 . Let the set of people be P and the set of people arranged on the same table as p_0 using our algorithm be A_0 . Let the optimal number of tables needed for a set of people be denoted by opt(.). Now, we will show that,

Claim 3.2.
$$opt(P) = opt(P \setminus A_0) + 1$$

Proof. We first note that from Claim 3.1, it is possible to generate a seating where no more swaps are possible. The set A_0 which is generated by our algorithm is one such arrangement for a single table. Thus, there exists an optimal arrangement with A_0 as one table. We will now prove the claim in two parts:

- 1. $opt(P) \le opt(P \setminus A_0) + 1$
 - An arrangement such that one table consists of people in A_0 and the remaining are arranged optimally will have the number of tables equal $opt(P \setminus A_0) + 1$. This is by definition \geq the most optimal seating for the set P and thus this equality trivially holds.
- 2. $opt(P \setminus A_0) \leq opt(P) 1$ We know that there exists an optimal arrangement which contains A_0 as one table. Therefore, the optimal arrangement for the remaining set of people, i.e.,

 $P \setminus A_0$ can be done using at least opt(P)-1 tables. Thus, we get that $opt(P \setminus A_0) \leq opt(P)-1$ which proves the second part of the claim.

From 1 and 2, we get that $opt(P) = opt(P \setminus A_0) + 1$, and this completes the proof of the claim.

Therefore, from Claim 3.2, we have shown that our algorithm greedily selects the most optimal seating. This completes the proof of correctness of our algorithm.