

CS213 2024 Tutorial Solutions

CS213/293 UG TAs

2024

Contents

Tutorial 1	2
Tutorial 2	7
Tutorial 3	12
Tutorial 4	17
Tutorial 5	23
Tutorial 6	28
Tutorial 7	31
Tutorial 8	35

Tutorial 1

1. The following is the code for insertion sort. Compute the exact worst-case running time of the code in terms of n and the cost of doing various machine operations.

Algorithm 1: Insertion Sort Algorithm

Data: Array A of length n

```
1 for  $j \leftarrow 1$  to  $n - 1$  do
2    $key \leftarrow A[j]$ ;
3    $i \leftarrow j - 1$ ;
4   while  $i \geq 0$  do
5     if  $A[i] > key$  then
6        $A[i + 1] \leftarrow A[i]$ ;
7     end
8     else
9       break;
10    end
11     $i \leftarrow i - 1$ ;
12  end
13   $A[i + 1] \leftarrow key$ ;
14 end
```

Solution: To compute the worst-case running time, we must decide the code flow leading to the worst case. When the if-condition evaluates to false, the control breaks out of the inner loop. The worst case would have the if-condition never evaluates false. This will happen when the input is in a strictly decreasing order.

Counting the operations for the outer iterator j (which goes from 1 to $n - 1$). Consider for each outer loop iteration:

1. Consider the outer loop without the inner loop. Then it has 1 Comparison (Condition in for loop), 1 Increment (1 Assignment + 1 Arithmetic) (Update in for loop), 3 Assignments (Line 2, 3, 11) and 2 Memory Accesses ($A[j]$ and $A[i+1]$) and 1 Jump. Hence $(C + 2Ar + 4As + 2M + J)$
2. while loop runs for j times (as i goes from $j - 1$ to 0). Without the if-else, the while loop still has to do 1 Comparison (while condition), 1 Decrement (1 Assignment + 1 Arithmetic) ($i \leftarrow i - 1$), 1 Jump, every iteration. Hence $(C + As + Ar + J) * j$
3. The if condition is checked in all iterations of the while loop. It involves 1 Comparison (if) and 1 Memory Access ($A[i]$ in the if condition). if block is executed for all iterations of while in the worst case. if block has 2 Memory Accesses and an Assignment and then a Jump happens to outside the if-else statement. Hence $(C + 3M + As + J) * j$

The total cost would then be (terms are for outer loop, inner loop, if and else):

$$\begin{aligned} & \sum_{j=1}^{n-1} ((C + 2Ar + 4As + 2M + J) + (C + As + Ar + J) * j + (C + 3M + As + J) * j) \\ &= (C + 2Ar + 4As + 2M + J) * (n - 1) + (C + As + Ar + J) * (n * (n - 1) / 2) \end{aligned}$$

$$+(C + 3M + As + J) * (n * (n - 1) / 2)$$

This means Insertion Sort is $\mathcal{O}(n^2)$.

(C is Comparison, As is Assignment, Ar is Arithmetic, M is Memory Access, J is Jump)

Note: Some intermediate operations might have been omitted but they do not affect the overall asymptotic performance of the algorithm. The time taken for comparisons, jumps, arithmetic operations, and memory accesses are assumed to be constants for a given machine and architecture.

2. What is the time complexity of binary addition and multiplication? How much time does it take to do unary addition?

Solution: Note that time complexity is measured as a function of the input size since we aim to determine how the time the algorithm takes scales with the input size.

Binary Addition

Assume two numbers A and B. In binary notation, their lengths (number of bits) are m and n. Then the time complexity of binary addition would be $\mathcal{O}(m + n)$. This is because we can start from the right end and add (keeping carry in mind) from right to left. Each bit requires an $\mathcal{O}(1)$ computation since there are only 8 combinations (2 each for bit 1, bit 2, and carry). Since the length of a number N in bits is $\log N$, the time complexity is $\mathcal{O}(\log A + \log B) = \mathcal{O}(\log(AB)) = \mathcal{O}(m + n)$.

Binary Multiplication

Similar to above, but here the difference would be that each bit of the larger number (assumed to be A without loss of generality) would need to be multiplied by the smaller number, and the result would need to be added. Each bit of the larger number would take $\mathcal{O}(n)$ computations, and then m such numbers would be added. So the time complexity would be $\mathcal{O}(m \times \mathcal{O}(n)) = \mathcal{O}(mn)^1$. Following the definition above, the time complexity is $\mathcal{O}(\log A \times \log B) = \mathcal{O}(mn)$.

Side note: How do we know if this is the most efficient algorithm? Turns out, this is NOT the most efficient algorithm. The most efficient algorithm (in terms of asymptotic time complexity) has a time complexity of $\mathcal{O}(n \log n)$ where n is the maximum number of bits in the 2 numbers.

Unary Addition

The unary addition of A and B is just their concatenation. This means the result would have A + B number of 1's. Iterating over the numbers linearly would give a time complexity of $\mathcal{O}(A + B)$ which is linear in input size.

3. Given $f(n) = a_0n^0 + \dots + a_dn^d$ and $g(n) = b_0n^0 + \dots + b_en^e$ with $d > e$, then show that $f(n) \notin \mathcal{O}(g(n))$

Solution: Let us begin by assuming the proposition is False, ergo, $f(n) \in \mathcal{O}(g(n))$. By definition, then, there exists a constant c such that there exists another constant n_0 such that

$$\forall n \geq n_0, f(n) \leq cg(n)$$

. Hence, we have

$$\forall n \geq n_0, a_0n^0 + \dots + a_dn^d \leq cb_0n^0 + \dots + b_en^e$$

$$\forall n \geq n_0, \sum_{i=0}^e (a_i - cb_i)n^i + a_{i+1}n^{i+1} + \dots + a_dn^d \leq 0$$

By definition of limit

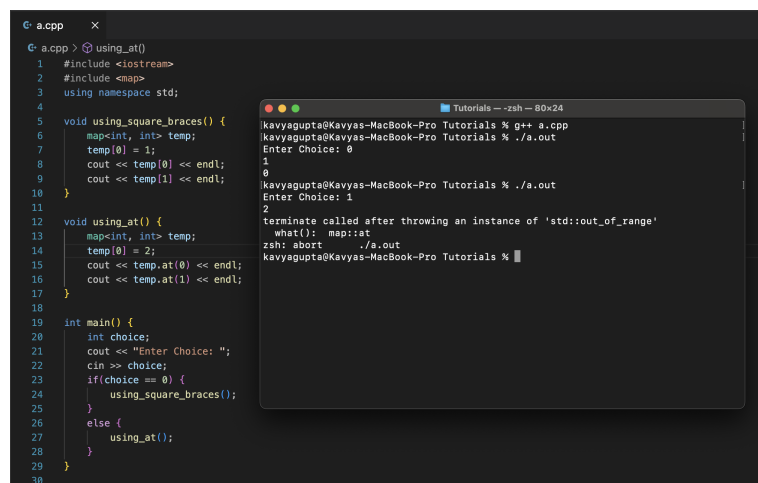
$$\lim_{n \rightarrow \infty} \sum_{i=0}^e (a_i - cb_i)n^i + a_{i+1}n^{i+1} + \dots + a_d n^d \leq 0 \\ \implies a_d \leq 0$$

Assuming $a_d > 0$ (since we are dealing with functions mapping from \mathbb{N} to \mathbb{N}), this results in a contradiction; thus, our original proposition is proved.

4. What is the difference between “at” and “[..]” accesses in C++ maps?

Solution: Both accesses will first search for the given key in the map but will behave differently when the key is not present. The at method will throw an exception if the key is not found in the map, but the [] operator will insert a new element with the key and a default-initialized value for the mapped type and will return that default value.

Look at the following illustration:-



```
g++ a.cpp
1 #include <iostream>
2 #include <map>
3 using namespace std;
4
5 void using_square_braces() {
6     map<int, int> temp;
7     temp[0] = 1;
8     cout << temp[0] << endl;
9     cout << temp[1] << endl;
10 }
11
12 void using_at() {
13     map<int, int> temp;
14     temp[0] = 2;
15     cout << temp.at(0) << endl;
16     cout << temp.at(1) << endl;
17 }
18
19 int main() {
20     int choice;
21     cout << "Enter Choice: ";
22     cin >> choice;
23     if(choice == 0) {
24         using_square_braces();
25     }
26     else {
27         using_at();
28     }
29 }
30
```

```
kavyagupta@Kavyas-MacBook-Pro:~/Tutorials$ g++ a.cpp
kavyagupta@Kavyas-MacBook-Pro:~/Tutorials$ ./a.out
Enter Choice: 0
1
0
kavyagupta@Kavyas-MacBook-Pro:~/Tutorials$ ./a.out
Enter Choice: 1
2
terminate called after throwing an instance of 'std::out_of_range'
what(): map::at
zsh: abort ./a.out
kavyagupta@Kavyas-MacBook-Pro:~/Tutorials$
```

5. C++ does not provide active memory management. However, smart pointers in C++ allow us the capability of a garbage collector. The smart pointer classes in C++ are

- auto_ptr
- unique_ptr
- shared_ptr
- weak_ptr

Write programs that illustrate the differences among the above smart pointers.

Solution: Memory allocated in heap (using new or malloc) if not de-allocated can lead to memory leaks. Smart pointers in C++ deal with this issue. Broadly, they are classes that store reference counts to the memory they point. When a smart pointer is created, reference count to that memory is increased by one. Whenever a smart pointer (pointing to the same memory) goes out of scope, its destructor is automatically executed, which reduces the reference count of that memory by one and when it hits zero, that memory is deallocated.

- `shared_ptr` allows multiple references to a memory location (or an object)
- `weak_ptr` allows one to refer to an object without having the reference counted
- `unique_ptr` is like `shared_ptr` but it does not allow a programmer to have two references to a memory location
- `auto_ptr` is just the deprecated version of `unique_ptr`

The use for `weak_ptr` is to avoid the circular dependency created when two or more object pointing to each other using `shared_ptr`.

You can see the reference count of a `shared_ptr` using its `use_count()` function.

```

C: shared_ptr.cpp X
1 #include <iostream>
2 #include <memory>
3 using namespace std;
4
5 int main() {
6     shared_ptr<int> P1(new int(4));
7     cout << P1.use_count() << endl;
8     shared_ptr<int> P2 = P1;
9     cout << P2.use_count() << endl;
10    cout << *P2 << endl;
11 }
12
Tutorials -- zsh -- 80x24
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ shared_ptr.cpp
kavyagupta@Kavyas-MacBook-Pro Tutorials % ./a.out
1
2
4
kavyagupta@Kavyas-MacBook-Pro Tutorials %

```

```

C: unique_ptr.cpp X
1 #include <iostream>
2 #include <memory>
3 using namespace std;
4
5 int main() {
6     unique_ptr<int> P1(new int(4));
7     unique_ptr<int> P2 = P1;
8 }
9
Tutorials -- zsh -- 80x24
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ unique_ptr.cpp
unique_ptr.cpp: In function 'int main()':
unique_ptr.cpp:7:26: error: use of deleted function 'std::unique_ptr<_Tp, _Dp>::unique_ptr(const std::unique_ptr<_Tp, _Dp>&) [with _Tp = int; _Dp = std::default_delete<int>]'
7 |     unique_ptr<int> P2 = P1;
  |     ~~~~~^~~~~
In file included from /opt/homebrew/Cellar/gcc/14.1.0_1/include/c++/14/memory:7
      from unique_ptr.cpp:2:
/opt/homebrew/Cellar/gcc/14.1.0_1/include/c++/14/bits/unique_ptr.h:516:7: note:
declared here
516 |     unique_ptr(const unique_ptr&) = delete;
    |     ~~~~~^~~~~
unique_ptr.cpp:7:26: note: use '-fdiagnostics-all-candidates' to display consid
red candidates
7 |     unique_ptr<int> P2 = P1;
  |     ~~~~~^~~~~
kavyagupta@Kavyas-MacBook-Pro Tutorials %

```

```

C: weak_ptr.cpp X
1 #include <iostream>
2 #include <memory>
3 using namespace std;
4
5 int main() {
6     shared_ptr<int> P1(new int(4));
7     cout << P1.use_count() << endl;
8     weak_ptr<int> P2 = P1;
9     cout << P1.use_count() << endl;
10 }
11
Tutorials -- zsh -- 80x24
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ weak_ptr.cpp
kavyagupta@Kavyas-MacBook-Pro Tutorials % ./a.out
1
1
kavyagupta@Kavyas-MacBook-Pro Tutorials %

```

```

C: cycle_error.cpp X
1 #include <iostream>
2 #include <memory>
3 using namespace std;
4
5 class B;
6
7 class A
8 {
9 public:
10     shared_ptr<B> sP1;
11     A() { cout << "A()" << endl; }
12     ~A() { cout << "~A()" << endl; }
13 };
14
15 class B
16 {
17 public:
18     shared_ptr<B> sP1;
19     B() { cout << "B()" << endl; }
20     ~B() { cout << "~B()" << endl; }
21 };
22
23 int main()
24 {
25     shared_ptr<A> aPtr(new A);
26     shared_ptr<B> bPtr(new B);
27
28     aPtr->sP1 = bPtr;
29     bPtr->sP1 = aPtr;
30
31     cout << aPtr.use_count() << endl;
32     cout << bPtr.use_count() << endl;
33 }
Tutorials -- zsh -- 80x24
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ cycle_error.cpp
kavyagupta@Kavyas-MacBook-Pro Tutorials % ./a.out
A()
B()
2
2
kavyagupta@Kavyas-MacBook-Pro Tutorials %

```

6. Why do the following three writes cause compilation errors in the C++20 compiler?

```
class Node {
public :
    Node(): value(o) {}
    const Node& foo(const Node* x) const {
        value = 3; // Not allowed because of -----
        x[0].value = 4; // Not allowed because of -----
        return x[0];
    }
    int value;
};

int main () {
    Node x[3], y ;
    auto &z = y.foo(x);
    z.value = 5; // Not allowed because of -----
}
```

Solution: The line 'value = 3;' is not allowed since the method 'foo' is marked **const** at the end. Using **const** after the function signature and its parameters for a class method implies that the class members cannot be changed and the object itself is constant within the bounds of the function. As a result, an error is raised.

The line 'x[0].value = 4;' raises an error since the parameter 'x' is of type '**const Node***', meaning that different values can be assigned to the pointer itself, but the subscripts of the pointer cannot be assigned, since it points to constant members of class 'Node'. Likewise, '*x.value = 3;' is equally illegal. Note that '**Node const ***' also does the same thing, whereas '**Node * const**' means that the subscripts of the pointer can be assigned since it does not point to constant members of class 'Node' but the pointer itself cannot be assigned to point to a different Node or array of Nodes. '**const Node* const**' or '**Node const * const**' imply that the pointer cannot be assigned AND subscripts or dereferences cannot be assigned.

The last line 'z.value = 5;' is illegal since the datatype of 'z', as determined from the method 'foo' is '**const Node &**' NOT '**Node &**'. The implication is that 'z' refers to a Node, which is constant and hence cannot be assigned. Note that assignments to class members are as bad as assignments to the class object itself regarding constant objects in C++.

Tutorial 2

1. The span of a stock's price on i^{th} day is the maximum number of consecutive days (up to i^{th} day and including the i^{th} day) the price of the stock has been less than or equal to its price on day i .

Example: for the price sequence 2 4 6 5 9 10 of a stock, the span of prices is 1 2 3 1 5 6.

Give a linear-time algorithm that computes s_i for a given price series.

Solution: The idea here is to find out the latest (most recent) day till the i^{th} day where the price was greater than the price on day i . To do this, we will maintain a stack of indices, and for each index i , we will keep the stack in a state such that the following invariant holds for each day i :

If the stack is not empty, then the price on the day whose index is at the top of the stack is the most recent price that was strictly greater than the current price.

With this invariant, the updates to the stack and the stock span follow as:

- we remove the top index in the stack till the price corresponding to the top index becomes strictly greater, or the stack becomes empty
- if the stack is empty, the current price is the largest, so the span is the number of days so far
- otherwise the span is the number of days since the index at the top of the stack

(Such a stack is known as a “Monotonic Stack”).

The base case is that at day 1 the price is the only encountered price, and hence is the largest, so the span is 1.

The time complexity is indeed $\mathcal{O}(n)$ (n being the number of days/size of array s) because each element is pushed **exactly once** in the stack and popped **atmost once**.

Look at Algorithm 2 for the pseudo-code.

Algorithm 2: Stock Span Algorithm

Data: Array *price* of length n and an empty stack *St*

```
1 for  $i \leftarrow 0$  to  $n - 1$  do
2    $s[i] \leftarrow i + 1$  // If  $\text{price}[i]$  is  $\geq$  all the previous prices, then  $s[i]$  will be  $i + 1$ 
3   while not St.isEmpty() do
4     if  $\text{price}[\text{St.top}()] > \text{price}[i]$  then
5        $s[i] \leftarrow i - \text{St.top}()$ 
6       break
7     end
8     St.pop()
9   end
10  St.push(i)
11 end
```

2. There is a stack of dosas on a tava, of distinct radii. We want to serve the dosas of increasing radii. Only two operations are allowed:

(i) serve the top dosa

(ii) insert a spatula (flat spoon) in the middle, say after the first k , hold up this partial stack and flip it upside-down and put it back

Design a data structure to represent the tava, input a given tava, and to produce an output in sorted order. What is the time complexity of your algorithm?

This is also related to the train-shunting problem.

Solution: Read the first line of the question with the emphasis being on **stack**. We will use the array representation of a stack for the tava.

Our stack-tava has the following abstraction:

- Initialization and input: Initialize an array S of size n , to represent the stack. Using the given input, fill the array elements with the radii of dosas on the tava in the initial stack order, with the bottom dosa at index 0 and the top at $n - 1$. Maintain a variable sz which contains the current size of the stack, and initialize it to n .
- Serving the top dosa: This method is essentially carrying out the “pop” operation on our stack S . Return the top dosa in the stack, $S[sz - 1]$, and decrement sz by 1. Clearly, this method takes a constant $\mathcal{O}(1)$ time for every call.
- Flipping the top k dosas: This method is equivalent to reversing the slice of the array S from index $sz - k$ to $sz - 1$. This is simple enough: initialize an index i to $sz - k$ and another index j to $sz - 1$, swap elements at indices i and j , increment i and decrement j by 1, and repeat the swap and increment/decrement while $j > i$. This method takes $\mathcal{O}(k)$ time for every call.

Note that in this implementation, the serve operation is the only way by which the problem can be reduced in size (number of dosas reduced by 1). The flip operation, if used wisely, can give us access to dosas that can help us reduce the problem size.

With this implementation, the algorithm to serve the dosas in increasing order of the radii can be designed as follows. While the stack is not empty, repeat the following:

- Iterate over the array S and find the index of the minimum element in the array. Let this index be m . This is the position of the dosa having the smallest radius
- Flip over the top $sz - m$ dosas in the stack, using the flipping operation. This brings the smallest dosa from index m to index $sz - 1$, id est, the top of the stack
- Serve the top dosa, using the serving operation. This pops the smallest dosa off the stack

It is easy to argue the correctness. The invariant is that we keep serving the minimum radii dosa from the remaining stack. The resulting order has to be sorted.

For computing the time complexity: in a stack of size z , finding the index of the minimum element takes $\mathcal{O}(z)$ time, flipping over the top $k \leq z$ elements takes $\mathcal{O}(k)$ time (and thus $\mathcal{O}(z)$ time) and serving off the top dosa takes $\mathcal{O}(1)$ time. Thus serving the smallest dosa takes $\mathcal{O}(z)$ time overall. This has to be repeated for all dosas, so the stack size z goes from n to 1, decrementing by 1 every time. Therefore the algorithm has a time complexity of $\mathcal{O}(n^2)$.

3. (a) Do the analysis of performance of exponential growth if the growth factor is three instead of two? Does it give us better or worse performance than doubling policy?
(b) Can we do the similar analysis for growth factor 1.5?

Solution: Let us do the analysis for a general growth factor α . Suppose initially $N = 1$ and there are $n = \alpha^i$ consecutive pushes. So, total cost of expansion is (refer to slides for detailed explanation):

$$(\alpha + 1) \cdot (\alpha^0 + \alpha^1 + \alpha^2 \dots + \alpha^{i-1})$$

$$\frac{\alpha + 1}{\alpha - 1} \cdot (\alpha^i - 1)$$

$$\frac{\alpha + 1}{\alpha - 1} \cdot (n - 1)$$

For α equals to 3 cost of expansion is $2 \cdot (n - 1)$, it's better than doubling policy. For α equals 1.5 cost of expansion is $5 \cdot (n - 1)$, which is worse than doubling policy. Trade-off involved is extra memory allocation, maximum extra memory allocated is $\alpha - 1$ times the requirement, so with increasing alpha extra memory allocation is increasing.

4. Give an algorithm to reverse a linked list. You must use only three extra pointers.

Solution: Let us initialise three pointers - *prev*, *curr*, and *next*, initialised to null, head and *head* \rightarrow *next* respectively. If head is null, then it is an empty linked list and we return. Otherwise, we will be iterating over each element in the linked list. In each iteration, perform the following updates.

- *curr* \rightarrow *next* = *prev*
- *prev* = *curr*
- *curr* = *next*
- *next* = *next* \rightarrow *next*

The loop terminates when *next* is null, in which case, set head to *curr*.

5. Give an algorithm to find the middle element of a singly linked list.

Solution: The idea is as follows:

- Make 2 pointers middle and end initialized to head of linked list
- At each step, increment middle by one and end by 2
- Continue this process until end reaches the end of linked list
- return the middle element
- return null if the head itself was null (empty linked list)

Note: When there is an odd number of elements, the above algorithm returns the element at the median location. What happens when there is an even number of elements?

6. Given two stacks S_1 and S_2 (working in the LIFO method) as black boxes, with the regular methods: "Push", "Pop", and "isEmpty", you need to implement a Queue (specifically: Enqueue and Dequeue working in the FIFO method). Assume there are n Enqueue/Dequeue operations on your queue. The

time complexity of a single method *Enqueue* or *Dequeue* may be linear in n , however the total time complexity of the n operations should also be $\Theta(n)$

Solution:

Let us try to formulate a simple implementation of a queue using 2 stacks.

Label the stacks as A and B . We will use stack A to store elements of the queue sequentially, and stack B to implement the dequeue operation.

Enqueue: To insert a new element to the queue, we can just *push* the element to the top of stack A . Here, each *enqueue* operation requires $\mathcal{O}(1)$ time.

Dequeue: To remove an element from the queue in a FIFO manner, we must somehow delete the first element pushed to stack A . Since we cannot access any element of the stack other than the top one (last inserted), we will have to use stack B to implement this.

To get access to the first inserted element, we need to *pop* all the remaining elements from stack A and store them into stack B . As you can see in the figure 1, the order of elements is reversed while moving them to stack B . When we reach the first element of A , we can just *pop* it without then adding it to B .

We must now decide what to do with the elements of stack B . One idea is to simply use stack B for storing the queue in the place of stack A , but this doesn't work since the order of elements would change.

Another possible approach would be to again transfer all the elements saved in stack B to stack A and restore the queue. This would preserve the order of elements and allow new elements to be pushed to A while maintaining correctness. However, think of the worst case time complexity of this approach. For every *dequeue* operation, all the elements have to be shifted twice - once from A to B and then from B to A . This causes each *dequeue* operation to be of $\mathcal{O}(n)$ time complexity and the total time complexity of n operations to be $\mathcal{O}(n^2)$. (think when?)

The blowup in time complexity is caused due to elements being shifted back to A from B on every *dequeue*. Do we really need that? What happens if we just let the elements stay in B ?

What happens now if we have 2 successive dequeue operations? Suppose we transferred the elements of queue from stack A to B in the first operation. For the 2nd dequeue operation, we can simply *pop* from stack B . Hence we don't need to re-transfer the elements. Now, if we have n operations, our complexity reduces to $\mathcal{O}(n)$.

To sum it up: Look at Algorithm 3 for the pseudo-code

Algorithm 3: Queue using Two Stacks

```
1 Function Enqueue(element):  
2   | S1.Push(element)  
3 end  
4 Function Dequeue():  
5   | if S1.isEmpty() and S2.isEmpty() then  
6     | throw("Queue is empty.")  
7   | end  
8   | if S2.isEmpty() then  
9     | while not S1.isEmpty() do  
10    |   | S2.Push(S1.Pop())  
11    | end  
12  | end  
13  | return S2.Pop()  
14 end
```

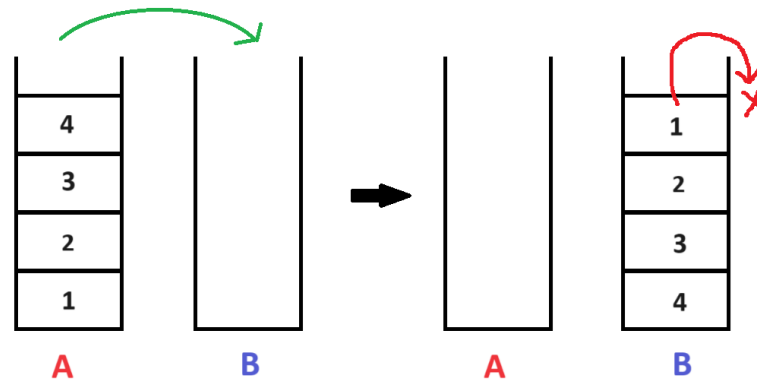


Figure 1: 2 Stacks making a Queue

Tutorial 3

1. What is the probability for the 3rd insertion to have exactly two collisions while using linear probing in the hash table?

Solution: Let the size of the hash table be $n \geq 3$. We assume that insertions take place via linear probing. Further, we assume that no deletions or insertions occur in the same sequence. Lastly, we assume that the hash function is uniform, i.e., for a random key, the likelihood of any hash value occurring is exactly $1/n$. After collisions, let the first two insertions occur at indices i and j .

If $((i + 1) \bmod n) \neq j$ and $((j + 1) \bmod n) \neq i$, then the third insertion cannot have two collisions. The reason is that if the hash value of the third key collides with the first key, then index $(i + 1) \bmod n$ is vacant, and there is only one collision. If the hash value is at index j , then index $(j + 1) \bmod n$ is vacant and hence, only one collision. Any other hash value implies no collision. Let us neglect this case.

We focus only on cases where i and j are adjacent to each other on the table. Let us re-label the indices as a and b such that $((a + 1) \bmod n) = b$. In this case, for any constant $0 \leq c < n$, the absolute value of $P(a = c)$ (not conditional) can be thought of as the probability that $((i = c \text{ and } j = (i + 1) \bmod n) \text{ or } (j = c \text{ and } i = (j + 1) \bmod n))$. The first event occurs when the second hash value is either i or $(i + 1) \bmod n$, with a chance of $(1/n) \cdot (2/n)$. The second event occurs when the second hash value is $(n + i - 1) \bmod n$, with a chance of $(1/n) \cdot (1/n)$.

From the above, we can say that $P(a = c) = (3/n^2)$ for all constants c , we safely say that two collisions occur if and only if the third hash-value is a , which occurs with a probability of $\frac{1}{n}$. However, c and, therefore, a can take any value. Thus, the product $(3/n^2) \cdot (1/n)$ has to be summed up across all values (n times), giving us the required answer: $(3/n^2)$.

2. Given that k elements have to be stored using a hash function with target space n . What is the probability of the hash function having an inherent collision? What is an estimate of the probability of a collision in the insertion of n elements?

Solution: To compute the probability, we first count the negative case: no collisions at all. This means that the k keys get mapped to distinct k keys out of n . We assume that for a random key and for any constant c , $P(\text{hash}(k) = c) = 1/n$.

The first key has n choices of hash values. Given each choice of the first key, the second key has $n - 1$ possible hash values such that there is no collision. Given the first two keys and hash values, the third has $n - 2$ choices, and so on, till the k^{th} key has $n - k + 1$ choices. Overall, the ordered k -tuple has $\prod_{i=0}^{k-1} (n - i)$ choices for no collision. All those are mutually exclusive events with likelihoods $(1/n)^k$ each, being the product of k probabilities of $1/n$ for k independent events and their respective outcomes. Thus,

$$P(\text{No collisions}) = \frac{n!}{(n - k)! \cdot n^k}$$

One minus this expression gives the likelihood of a collision.

If $k = n$, the expression boils down to $(n!/n^n)$. We apply Stirling's approximation here:

$$\sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n+1}} < n! < \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n}}$$

$$\sqrt{2\pi n} \left(\frac{1}{e}\right)^n e^{\frac{1}{12n+1}} < \frac{n!}{n^n} < \sqrt{2\pi n} \left(\frac{1}{e}\right)^n e^{\frac{1}{12n}}$$

For large n , we can neglect the terms $e^{(1/(12n+1))}$ and $e^{(1/12n)}$ as both are very close to 1. However, the term $\sqrt{2\pi n} \cdot e^{-n}$ tends to zero as n grows very large. This means that the required probability is sandwiched between terms that exponentially decay toward zero.

The implication of the above is the likelihood that all n keys get uniquely mapped to the n possible hash values decay exponentially with n , and thus, collisions are almost certain for large n . This can be linked with the birthday paradox in that in a class of 365; it is almost impossible for all students to have different birthdays; even a single tutorial batch has a high likelihood of at least two birthdays being the same. In the context of hashing, although the expected number of terms sharing a hash value is one per hash value, we expect a good number of collisions to occur during hashing.

3. Let $C(i)$ be the chain of array indices that are queried to look for a key k in linear probing where $h(k) = i$.
 - (a) How does this chain extend by an insertion, and how does it change by a deletion?
 - (b) A search for a key k ends when an empty cell is encountered. What if we mark the end of $C(i)$ with an end marker? We stop the search when this marker is encountered. Would this work? Would this be efficient?
 - (c) Is there a way of not using tombstones?

Solution: To recall, linear probing covers indices $i, i + 1, \dots$ until an empty index is found (then k not found in the table) or the cell stores k . Note that tombstones do not count as empty cells in the hash table; counting them so after deletion would break existing chains and result in false negatives when searching. Including them means no change in these chains, answering part of the first (a) subproblem.

If the chain terminates at k (key found), then inserting a new element does not change anything. If the added element is k itself, then it occupies the position at the end of the chain, a vacant cell. So, no change. In the other case, the chain extends by one if the inserted element is at the end of the chain and is not k . The chain could extend a lot more if the inserted element joins the existing chain with another series of filled cells in the table. Note that the chain, except for the end, has only filled or tombstone cells; thus, the insertion cannot occur in the middle of the chain. If the insertion of the element that is not k does not occur at the end of the previous chain for searching k , then there is no change.

The answer to the second (b) part is that it would work but not be any more efficient than the standard algorithm. For this to work, an empty cell once encountered during the search should be marked an 'end.' Otherwise, the first search would go on an infinite loop, looking for an end marker. If this is the case, the search must look for both end marks and empty cells in each iteration, making it slightly less efficient but at the same time complexity.

For the third (c) part, a slight modification should be made to the algorithm discussed above: all cells should be marked 'end' at the start. Otherwise, a deletion without a tombstone would cause false negatives in future searches, just like the discussions without tombstones. If this modification is made, then we basically replace "empty" with an "end" marker and "tombstone" marker with "empty". The algorithm is, otherwise, identical to the normal algorithm with tombstones, and hence, we really do

not avoid tombstones, although we do not use them. Note that marking the freed cells as 'end' on deletions would be as bad as not using tombstones, which results in false negatives.

4. Let $m = 11$, $h_1(k) = (k \bmod 11)$, $h_2(k) = 6 - (k \bmod 6)$.
Let us use the following hash function for an open addressing scheme.

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$$

1. What will be the state of the table after insertions of 41, 22, 44, 59, 32, 31, and 74?
2. Let $h_2(k) = p - (k \bmod p)$. What should be the relationship between p and m such that h is a valid function for linear probing?
3. What is the average number of probes for an unsuccessful search if the table has α load factor?
4. What is the average time for a successful search?

Solution: To recall, the standard linear probing uses the function $h(k, i) = (h(k) + i) \bmod m$ where m is the size of the hash table [Here 11]. The way it works is that we check the cells indexed at $h(k, 0), h(k, 1), h(k, 2), \dots$ till we find an empty cell $h(k, a)$ for some a . The same logic should be extended here, except $h(k, i)$, which is defined in the problem statement.

Let the table before adding 41 be

--	--	--	--	--	--	--	--	--	--	--

Now, $h_1(41) = 8$ and $h_2(41) = 1$. This means that $h(41, 0) = 8$ and cell 8 is empty. So, we add 41 there:

								41		
--	--	--	--	--	--	--	--	----	--	--

$h_1(22) = 0$ and cell 0 is empty so we can add 22 there. [$h(k, 0) = h_1(k)$]

22								41		
----	--	--	--	--	--	--	--	----	--	--

$h_1(44) = 0$ but cell 0 is filled. Now, $h_2(44) = 4$ and $h(44, 1) = 4$ which is an empty tile.

22				44				41		
----	--	--	--	----	--	--	--	----	--	--

$h_1(59) = 4$ which is a filled cell. $h_2(59) = 1$ and $h(59, 1) = 5$ is an empty tile.

22				44	59			41		
----	--	--	--	----	----	--	--	----	--	--

$h_1(32) = 10$ and $h_1(31) = 9$ are both empty and different cells, so we can add them.

22				44	59			41	31	32
----	--	--	--	----	----	--	--	----	----	----

Lastly, $h_1(74) = 8$, which is occupied. $h_2(74) = 4$. Now, $h(74, 1) = 1$ which is free.

22	74			44	59			41	31	32
----	----	--	--	----	----	--	--	----	----	----

This is how the table looks after inserting the elements in the given order.

Moving on to the next part, a valid function for linear probing should eventually cover all cells in the grid as i increments. This is not possible if $\gcd(h_2(k), m) > 1$, for one or more values of k . The reason is that $h(k, i) \bmod \gcd(h_2(k), m)$ will be fixed, and hence, $h(k, i) \bmod m$ will only take some values. But if that gcd is 1, then it means that for all $0 \leq i < m$, $\gcd \cdot i \bmod m$ will have a unique value, and the number of these unique values will be m , so all cells of the table will be covered.

Note that $\gcd(h_2(k), m) = 1$ should hold for all values of k , not just for any one k . Keeping $h_2(k) = p - (k \bmod p)$ can take any value in $1, 2, \dots, p$ depending on the value of k . Thus, none of the first p natural numbers should share a common factor with m . In other words, the smallest prime factor (not 1) of m , the size of the hash table, should be larger than p . It is not necessary that m is prime – for $p = 6$, $m = 7 \cdot 13 = 91$ also works.

Exact expression of expected number of probes in a search (unsuccessful/successful) is mathematically complex hence we'll prove an appropriate upper bound for them. Claim: The mentioned expression for an unsuccessful one is atmost $1/(1 - \alpha)$. This is assuming $\alpha < 1$.

Proof: Let X be the number of probes made in an unsuccessful search, and let A_i be the event that an i^{th} probe occurs and it is to an occupied slot. Then $X \geq i$ when the event $A_1 \cap A_2 \cap \dots \cap A_{i-1}$ occurs. Hence,

$$P(X \geq i) = P(A_1 \cap A_2 \cap \dots \cap A_{i-1}) \\ = P(A_1) \times P(A_2|A_1) \times P(A_3|A_2 \cap A_1) \times \dots \times P(A_{i-1}|A_1 \cap \dots \cap A_{i-2})$$

$P(A_1) = n/m$ because there are n elements and m slots. For $j > 1$, the probability that there is a j^{th} probe and it is to an occupied slot, given that the first $j - 1$ probes were to occupied slots, is $(n - j + 1)/(m - j + 1)$ because we are finding one of the remaining $(n - (j - 1))$ elements in one of the $(m - (j - 1))$ unexamined slots, and we have uniform hashing so each of the slots is equally likely to be chosen next.

Because we have $n < m$, we have $(n - j)/(m - j) \leq n/m \forall 0 \leq j < m$. Thus,

$$P(X \geq i) = \frac{n}{m} \times \frac{n-1}{m-1} \times \dots \times \frac{n-i+2}{m-i+2} \leq \left(\frac{n}{m}\right)^{i-1} = \alpha^{i-1}$$

So $E[X] = \sum_{i=1}^{m-1} P(i^{th} \text{ probe occurs}) = \sum_{i=1}^{m-1} P(X \geq i) \leq \sum_{i=1}^{m-1} \alpha^{i-1} \leq \sum_{i=1}^{\infty} \alpha^{i-1} = 1/(1 - \alpha)$. Hence proved.

Claim: Expected number of probes for a successful search is atmost $\frac{1}{\alpha} \ln(\frac{1}{1-\alpha})$.

Proof: A search for a key k reproduces the same probe sequence as when the element with key k was inserted. If k was the $(i + 1)^{th}$ key inserted into the hash table, there were only i keys in the table at the time, so the expected number of probes made in a search for k is at most $1/(1 - i/m) = m/(m - i)$ (as an insertion is equivalent to an unsuccessful search, here $\alpha = i/m$). Averaging over all n keys in the hash table, we get,

$$\begin{aligned} \frac{1}{n} \sum_{i=0}^{n-1} \frac{m}{m-i} &= \frac{m}{n} \sum_{i=0}^{n-1} \frac{1}{m-i} \\ &= \frac{1}{\alpha} \sum_{k=m-n+1}^m \frac{1}{k} \\ &\leq \frac{1}{\alpha} \int_{m-n}^n \frac{1}{x} dx \\ &= \frac{1}{\alpha} \ln\left(\frac{m}{m-n}\right) \\ &= \frac{1}{\alpha} \ln\left(\frac{1}{1-\alpha}\right) \end{aligned}$$

Hence Proved.

5. Suppose you want to store a large set of key-value pairs, for example, $(name, address)$. You have operations in this set: addition, deletion, and search of elements. You also have queries about whether a particular name or address is in the set, and if so, then count them and delete all such entries. How would you design your hash tables?

Solution: Let us recall that in the classic hashing, the hash function operates on the key, and the value is simply stored alongside the key. An intuitive idea is to keep the key (from the point of view of hashing) as **both** the name and the address, with no special value otherwise. The hash function operates on this and gives us the index. Unfortunately, this gives rise to a problem: there are multiple different hash values for a given name, and hence, searching for the existence of an entry given a name (but not address – we do not know which addresses are there given this name) is as bad as linearly iterating through the table.

One way to get around this is to use a two-level hash table: the outer table N stores the keys as the names alone, and each value corresponding to a name points to a hash table itself. This inner hash table $N[name]$ stores only addresses and only addresses where the name is the corresponding key of the outer table. This makes finding elements based on names easy but does not solve the problem of searching for given addresses.

To solve that, we can use a double two-level hash table: one two-level table N as described above and one two-level table A done the other way round, i.e., A stores the keys as addresses and values as pointers to inner hash tables $A[address]$ that store names corresponding to addresses. Note that this requires additional storage space (double) and more maintenance.

Specifically, it is important to ensure integrity: for every $(name, address)$ pair in the set, $N[name]$ should exist and contain $address$, and $A[address]$ should exist and contain $name$. Insertion should be done in both tables at once. On the other hand, if any $(name, address)$ pair does not exist, then $A[address]$, if it exists, should not contain $name$, and $N[name]$, if it exists, should not contain $address$. This should be reflected in the deletion of pairs.

Tutorial 4

1. Given a tree with a maximum number of children as k . We give a label between 0 and $k - 1$ to each node with the following simple rules.

(i) The root is labelled 0.

(ii) For any vertex v , suppose that it has r children, then arbitrarily label the children as $0, \dots, r - 1$.

This completes the labelling. For such a labelled tree T , and a vertex v , let $Seq(v)$ be the labels of the vertices of the path from the root to v . Let $Seq(T) = \{Seq(v) \mid v \in T\}$ be the set of label sequences.

What properties does $Seq(T)$ have? If a word w appears, what words are guaranteed to appear in $Seq(T)$? How many times does a word w appear as a prefix of some words in $Seq(T)$?

Solution: Some properties are $Seq(T)$ are:-

1. Size of this set will be n (n is the number of nodes in the tree) (See that for different v_1, v_2 vertices, $Seq(v_1) \neq Seq(v_2)$).
2. If there are h_j nodes of j length in set $Seq(T)$, then the tree has h_j nodes at j^{th} level.
3. If we consider ordering of set $Seq(T)$ in lexicographical order as the metric, then it will just represent a DFS on the tree.
4. If we consider ordering of set $Seq(T)$ by length followed by lexicographical order as the metric, then it will just represent a BFS on the tree.

$p_0 = 0$ (root of tree)

$$w = p_0 p_1 p_2 \dots p_n$$

$$S_w = \{x \mid x = p_0 p_1 p_2 \dots p_i q \text{ where } 0 \leq i \leq n - 1 \text{ \& } 0 \leq q \leq p_{i+1}\} \cup \{0\}$$

We can prove S_w is the set which is guaranteed to be present if w is present by following argument:- Each element of set S_w is guaranteed to be present if w is present. We can prove this by considering two claims,

1. Any prefix of w will be present since this just represents the path from root to ancestor of w .
2. If a word $l_0 l_1 l_2 \dots l_n$ is present then all the words of form $l_0 l_1 l_2 \dots l_{n-1} q$ where $0 \leq q \leq l_n$ will also be present since these are just the child of the node represented by sequence $l_0 l_1 l_2 \dots l_{n-1}$

If $w = Seq(v)$ for some $v \in T$ and the number of nodes in the subtree of v be x , then w appears as a prefix for x number of words in $Seq(T)$.

2. Write a function that returns $lca(v, w, T)$. What is the time complexity of the program?

Solution: This is for a general tree and not just a binary tree.

See the Pseudo-code in Algorithm 4.

At any node (say x), if you find the LCA in x 's subtree, return it. Otherwise, see if v, w are present in x 's subtree (subtree includes x). If yes, then x is the LCA, otherwise return $NULL$.

Time complexity is $\mathcal{O}(n)$, n being the number of nodes of T , as we go through each node atmost once.

If there existed a "parent" pointer in each node, then can you form a new LCA algorithm ?

Algorithm 4: LCA

```

1 Function LCA( $curr, v, w$ ):
2   if  $curr = NULL$  then
3     | return  $False, False, NULL$  //This return type is an ordered tuple
4   end
5    $v\_present \leftarrow (curr = v)$  // True or False
6    $w\_present \leftarrow (curr = w)$ 
7   for  $child \in children[curr]$  do
8     |  $v\_in\_child, w\_in\_child, LCA\_child \leftarrow LCA(child, v, w)$ 
9     | if  $LCA\_child \neq NULL$  then
10      | return  $True, True, LCA\_child$ 
11    | end
12     $v\_present \leftarrow v\_present \parallel v\_in\_child$ 
13     $w\_present \leftarrow w\_present \parallel w\_in\_child$ 
14  end
15  if  $v\_present = True$  and  $w\_present = True$  then
16    | return  $True, True, curr$ 
17  end
18  return  $v\_present, w\_present, NULL$ 
19 end
20 // Call LCA( $root, v, w$ ) in the main function

```

3. Given $n \in T$, Let $f(n)$ be a vector, where $f(n)[i]$ is the number of nodes at depth i from n .
- Give a recursive equation for $f(n)$.
 - Give a pseudo-code to compute the vector $f(\text{root}(T))$. What is the time complexity of the program?

Solution:

Observations:

- For any node n , the only node at depth 0 from n is the node n itself.
- For a leaf node l , there is no node n such that $\text{depth}(n)$ from l is greater than or equal to 1.
- For a node n , if node m is at depth k , then there exists a *unique path* from n to m of length k , and this path must go through one of the children of n (let's say p) i.e. node m is at a depth $k - 1$ from node p .
Conversely, if a node m is at a depth $k - 1$ from node p , then it must be at a depth of k from node n where n is parent of p .

With these observations in mind, we can say that $f(n)[i] = \sum_{m \in \text{children}(n)} f(m)[i - 1]$. Recursive equation of $f(n)$ then can be written as:

$$f(n) = [1] + \left(\sum_{m \in \text{children}(n)} f(m) \right)$$

which is nothing but concatenating $[1]$ in front of the sum obtained from the f vectors of the children.

Here base case will be when n is a leaf, for that we have

$$f(n) = [1, 0, 0, \dots, 0]$$

Pseudocode: (Look at Algorithm 5)

As you all know, we will be doing some kind of traversal to get $f(\text{root}(T))$!!!

Can you think of the traversal method we gonna use by looking at the recursion obtained in **part 1** ???

I hope you got it correct, we are going to use post-order traversal (why so?)

We will use array data structure for storing $f(i)$ of size n , where n is the number of nodes in the tree

The algorithm given above, does a post order traversal on tree (note that this algorithm holds good for any tree and not just a binary tree), on getting to a leaf, it returns an array $[1, 0, \dots, 0]$, and on an internal node, it returns f according to the recursion obtained above.

Time complexity

Assume that vectors at all the nodes are of size n . There will be $\mathcal{O}(n)$ edges and for each edge we will do addition of two such vectors. Hence this algorithm has $\mathcal{O}(n^2)$ time complexity.

Algorithm 5: Computing $f(n)$

Data: Tree T having n nodes

```

1 Function ModifiedPostOrderWalk( $p$ ):
2    $f \leftarrow [1, 0, \dots, 0]$ 
3   if  $p.\text{isLeaf}()$  then
4     return  $f$ 
5   end
6   for  $n' \in \text{children}[p]$  do
7      $f_c \leftarrow \text{ModifiedPostOrderWalk}(n')$ 
8     for  $j \leftarrow 1$  to  $n - 1$  do
9        $f[j] \leftarrow f[j] + f_c[j - 1]$ 
10    end
11  end
12  return  $f$ 
13 end
```

4. Give an algorithm for reconstructing a binary tree if we have the preorder and inorder walks.

Solution: The key idea here is to note that we will recursively use both traversals in parallel. We will use the inorder traversal for finding the node index and iterate over the preorder to build the tree.

We assume that we have a function that can search for a node in the inorder traversal given the start and end indices of the search range.

To utilise both the traversals, we maintain a start and end index for the inorder traversal of the tree and make recursive calls while moving these indices and the iterator of the preorder traversal.

The tree build function can be written recursively as follows:

- if the start index is greater than the end index the tree is *NULL*
- otherwise we first create the root as the first element of the preorder
- if the start index and the end index are equal then the root is the answer
- otherwise we search for the index of the root in the inorder traversal
- the left sub-tree can be built recursively by setting the end index as the index previous to the root while the right sub-tree can be built after this by setting the start index as the index next to the root (note that the preorder traversal will be linearly iterated over by both these recursive calls but the iterator will be common across the calls)
- we return the root after setting the results of the sub-tree computations

Explanation: Note that if we know the position of the root in the inorder traversal, the left side of the inorder traversal is the inorder traversal of the left sub-tree and similarly for the right sub-tree. We also note that the first element in the preorder traversal is always the root, and the following elements will be the preorder traversal of the left sub-tree followed by the right sub-tree. The time complexity of this will be $\mathcal{O}(n^2)$ (think what is the worst case).

Correctness: Arguing for correctness is straightforward. Within each recursive call, the root is the element pointed to by the preorder iterator, which increments by one for each recursive call. After that, we find the root in the inorder traversal and then build the tree using the left side of the inorder traversal and the next set of the elements in the preorder traversal (both of these correspond to the left sub-tree). Then we build the tree using the right side of the inorder traversal and the remaining set of the elements in the preorder traversal (both correspond to the right sub-tree). Since the invariants (left-root-right for inorder and root-left-right for preorder) are maintained, the recursion terminates correctly.

Additional: What is the **minimal** set of changes required for building a tree using the postorder and inorder traversals?

5. Let us suppose all internal nodes of a binary tree have two children. Give an algorithm for reconstructing the binary tree if we have the preorder and postorder walks.

Solution:

Let preorder walk be $f = (f_1, f_2, \dots, f_n)$ and postorder walk be $l = (l_1, l_2, \dots, l_n)$

Base case

f and l have just 1 element let's say e , then the tree is nothing but a single-node tree with root as e .

Recursion step

Observations for tree with each internal node having two children and number of nodes greater than 1

1. In preorder walk, *first* element is *root* and *second* element is left child of root.

2. In postorder walk, *last* element is *root* and *second last* element is right child of root.

So, now with these observations in mind let's try to build the tree from f and l .

We have the root as $r = f_1 = l_n$

Left child of root as f_2 and right child of root as l_{n-1}

Now we will break f and l so as to get preorder and postorder walks for left and right subtree

Given f , which element in right subtree of r is first to be in the f ?

By Observation 1, preorder on right subtree (which is a tree) of r should have its root i.e. right child of r which is l_{n-1} first in the walk and so we have the break of f !!!

Find l_{n-1} in f , let us say f_k , then we have $f' = (f_2, f_3, \dots, f_{k-1})$ as preorder walk of left subtree of r and $f'' = (f_k, f_{k+1}, \dots, f_n)$ as preorder walk of right subtree of r

Similarly, given l , which element in left subtree of r is the last to be in the l ?

By Observation 2, postorder on left subtree (which is a tree) of r should have its root i.e. left child of r which is f_1 last in the walk and so we have the break of l !!!

Find f_2 in l , let us say l_k , then we have $l' = (l_1, l_2, \dots, l_k)$ as postorder walk of left subtree of r and $l'' = (l_{k+1}, l_{k+2}, \dots, l_{n-1})$ as postorder walk of right subtree of r

With f' and l' , we **recursively** construct the left subtree and with f'' and l'' , the right subtree. The time complexity will again be $\mathcal{O}(n^2)$.

6. (a) Show that in order printing of BST nodes produces a sorted sequence of keys.
 (b) Give a sorting procedure using BST.
 (c) Give the complexity of the procedure.

Solution:

- (a) At any non-leaf BST node, the key at any node in the left subtree is not greater than the key at the node. Similarly, the key at any node in the right subtree is not less than the key at the node. An inorder traversal of a BST visits the left subtree first, then the node, and then the right subtree. Thus, the inorder traversal of the BST will produce a sorted sequence of keys.
- (b)
- Create a BST with the first element of the unsorted array as the root and no children.
 - Insert each of the remaining elements into the BST.
 - Perform an inorder traversal of the BST to get the sorted sequence of keys.
- (c) The complexity of the procedure is $\mathcal{O}(n \log n)$. This is because average insertion of n elements in a BST takes $\mathcal{O}(n \log n)$ time. And the inorder traversal takes $\mathcal{O}(n)$ time.

7. Given a BST tree T and a value v , write a program to locate the leftmost and rightmost occurrence of the value v .

Solution: Leftmost occurrence in a BST means leftmost occurrence in the inorder traversal of that BST.

Look at Algorithm 6 for the Pseudo-code of the leftmost occurrence in a BST.

The code itself is self-explanatory. Similarly you can find this for the rightmost occurrence.

Time complexity is $\mathcal{O}(h)$, h being the height of T , as we go through each level atmost once.

Algorithm 6: Finding the leftmost occurrence

```
Data: Tree  $T$  having  $n$  nodes
1 Function Leftmost( $curr, v$ ):
2   if  $curr = NULL$  then
3     return  $NULL$ 
4   end
5   if  $curr \rightarrow val \geq v$  then
6      $left\_check \leftarrow Leftmost(curr \rightarrow left, v)$ 
7     if  $left\_check \neq NULL$  then
8       return  $left\_check$ 
9     end
10  end
11  if  $curr \rightarrow val = v$  then
12    return  $curr$ 
13  end
14  if  $curr \rightarrow val \leq v$  then
15     $right\_check \leftarrow Leftmost(curr \rightarrow right, v)$ 
16    if  $right\_check \neq NULL$  then
17      return  $right\_check$ 
18    end
19  end
20  return  $NULL$ 
21 end
22 // Call  $Leftmost(root, v)$  in the main function
```

Tutorial 5

1. Given a BST T and a key k , the task is to delete all keys $< k$ from T .

- (a) Write pseudocode to do this.
- (b) How much time does your algorithm take?
- (c) What is the structure of the tree left behind?
- (d) What is its root?

Solution:

- (a) The code (See Algorithm 7) provided deletes all keys $< k$ from T and returns the root of the modified BST.
- (b) The complexity of the procedure is $\mathcal{O}(h)$, h being the height of T , as we are visiting each level at most once (worst case being reaching a leaf node).
- (c) The structure of the tree left behind is still a BST.
- (d) If the previous root was greater than or equal to k , then the new root is the same as the previous root. Otherwise, the new root is different.

Algorithm 7: Deleting all keys $< k$

```
1 Function DeleteKeys(curr, k):  
2   if curr = NULL then  
3     return NULL  
4   end  
5   if curr → val ≥ k then  
6     (curr → left) ← DeleteKeys((curr → left), k)  
7     return curr  
8   end  
9   else  
10    return DeleteKeys((curr → right), k)  
11  end  
12 end  
13 // Call DeleteKeys(root, k) in the main function
```

2. Let $H(n)$ be the expected height of the tree obtained by inserting a random permutation of $[n]$. Write the recurrence relation for $H(n)$.

Solution: The height of a binary tree equals one plus maximum of that of left and right subtrees. In a random permutation of $[n]$, any element (say k) can be the root of the subtree with probability $1/n$. The left subtree will contain $k - 1$ nodes and the right one $n - k$.

Let X_n be the random variable representing the height of a n -node BST. Also, $H(n) = E[X_n]$. Then the recurrence relation is given as :-

$$\begin{aligned}
\forall n \in \mathbb{N}, H(n) &= E[X_n] = 1 + \frac{1}{n} \left(\sum_{i=1}^n E[\max(X_{i-1}, X_{n-i})] \right) \\
&\leq 1 + \frac{1}{n} \left(\sum_{i=1}^n E[X_{i-1} + X_{n-i}] \right) \\
&= 1 + \frac{1}{n} \left(\sum_{i=1}^n E[X_{i-1}] + \sum_{i=1}^n E[X_{n-i}] \right) \\
&= 1 + \frac{2}{n} \left(\sum_{i=1}^{n-1} E[X_i] \right) = 1 + \frac{2}{n} \left(\sum_{i=1}^{n-1} H_i \right)
\end{aligned}$$

Base Case: $H(0) = 0$

Solving the recurrence gives $H(n)$ as $\mathcal{O}(\log n)$. You can refer [this](#) for a rough idea of its solution.

3. Prove that after rotation the resulting tree is a binary search tree.

Solution: See Figure 2. Suppose we go from left to right.

A becomes the left child of B and T_2 becomes the right child of A . $A \rightarrow val \leq B \rightarrow val$ as B was initially the right child of A (definition of BST) and $T_2 \rightarrow val \geq A \rightarrow val$ as T_2 was originally in the right subtree of A . T_1 is a BST, so is T_2 , hence A 's subtree is a BST too. Hence the new structure follows the rules of a BST.

Similarly we can prove for right to left in the figure.

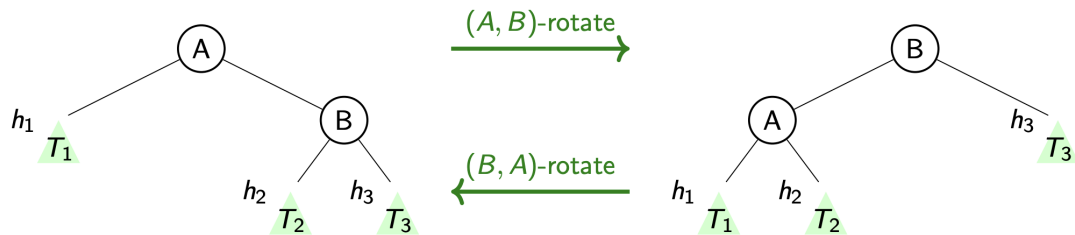


Figure 2: Rotation

4. Insert sorted numbers $1, 2, 3, \dots, 10$ to an empty red-black tree. Show all intermediate red-black trees.

Solution: See Figure 3 below.

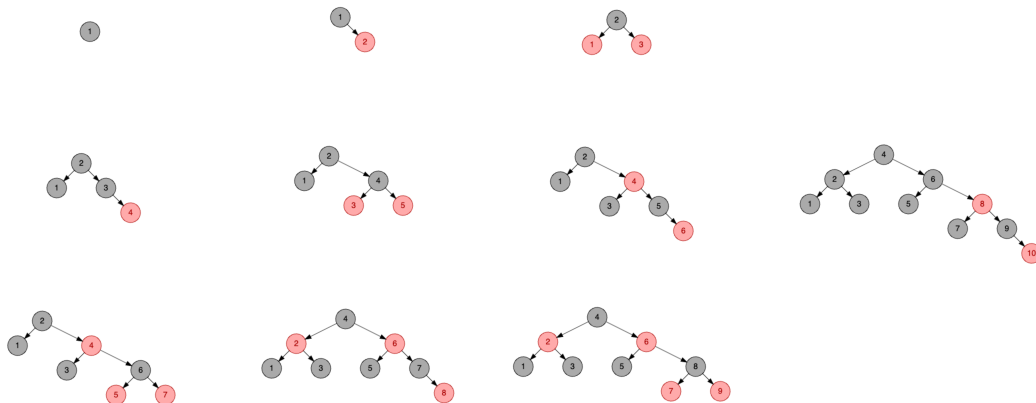
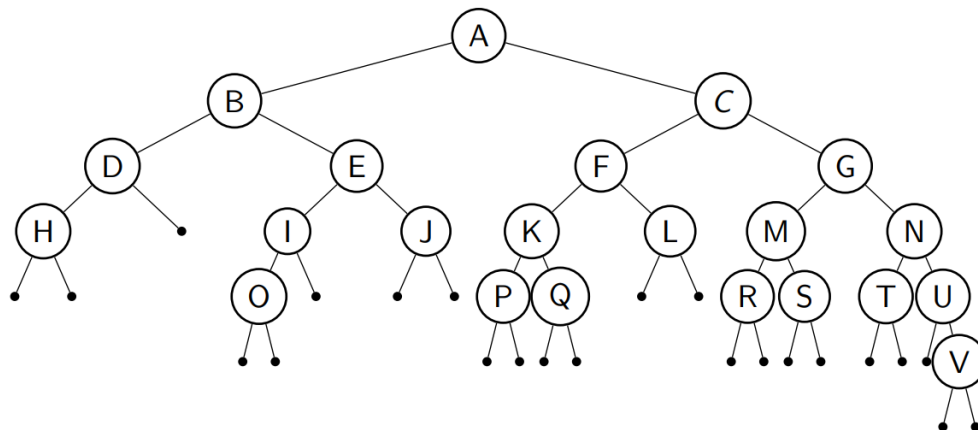


Figure 3: RBT Tree Insertions

5. Consider the tree below. Can it be coloured and turned into a red-black tree? If we wish to store the set $1, \dots, 22$, label each node with the correct number. Now add 23 to the set and then delete 1. Also, do the same in the reverse order. Are the answers the same? When will the answers be the same?



Solution: If we color this RB tree, Black height of this tree must be 3. Why?? (Hint: Argue in terms of property of RB tree) Colour the nodes H, E, O, C, N, V, K and M red; the rest should be black. This is a red-black tree with a black height of 3. Using the inorder traversal of the tree,

$A = 8, B = 3, C = 14, D = 2, E = 6, F = 12, G = 18, H = 1, I = 5, J = 7, K = 10$

$L = 13, M = 16, N = 20, O = 4, P = 9, Q = 11, R = 15, S = 17, T = 19, U = 21, V = 22$

Insertion Followed by Deletion

The insertion is very straightforward as 23 gets added as the right child of $V = 22$, let's call it W. Red black tree condition gets violated as V and W are of same color now. So, according to case 3, apply a left rotation such that after rotation $N.\text{right} = V$, $V.\text{left} = U$ and $V.\text{right} = W$. $H=1$ is a red color leaf, so we can simply delete it without any trouble of black height violation.

Deletion Followed by Insertion

In this case, our answer will be the same as the one above. Since, deletion of H doesn't change anything significant for the right tree and we can simply insert 23 as described above.

Difference can be in cases where there is upward propagation of violation while dealing

with violation after insertion and deletion. Since, upward propagation checks and affect different parts and colors of the tree

6. (a) Give running time complexities of delete, insert, and search in red-black tree.
 (b) What are the advantages of red-black tree as compare to Hash table, where every thing is constant time?

Solution: (a) $\mathcal{O}(\log n)$, n being the number of nodes of RBT. (Refer slides to know why).

(b) These are the following advantages:-

- Red-black trees maintain elements in a sorted order, which is useful for operations that require order (e.g., finding the minimum/maximum, range queries). Hash tables, on the other hand, do not maintain any order among elements.
- The performance of a red-black tree is predictable and does not depend some factor like hash table does on the quality of the hash function or the distribution of data.
- RBT has no collisions like a Hash Table.

7. Suppose we use a hash function h to hash n distinct keys into an array T of length m . Assuming simple uniform hashing, what is the expected number of collisions? More precisely, what is the expected cardinality of $\{\{k, l\} : k \neq l \text{ and } h(k) = h(l)\}$?

Solution: Let's define a binary random variable X_{ij} :-

$$X_{ij} = \begin{cases} 1 & \text{if there is a collision between the } i^{th} \text{ and } j^{th} \text{ keys during insertion,} \\ 0 & \text{otherwise.} \end{cases}$$

$E[X_{ij}] = 1/m$ as the probability that the key that is inserted later gets the same slot as the one inserted earlier is $1/m$. Our hashing function is **uniform**, i.e. probability of getting any slot is equal and hence is $1/(\text{number of slots})$.

Let X be the random variable returning the number of collisions during insertion of the n keys, then,

$$X = \sum_{i=1}^n \sum_{j=i+1}^n X_{ij}$$

This is because the RHS represents the sum over all possible distinct unordered pairs $\{\{k, l\}, k \neq l\}$ seeing if there is a collision within each pair or not. Hence,

$$\begin{aligned} E[X] &= E\left[\sum_{i=1}^n \sum_{j=i+1}^n X_{ij}\right] \\ &= \sum_{i=1}^n \sum_{j=i+1}^n E[X_{ij}] \text{ (Due to linearity of the Expectation operator)} \\ &= \sum_{i=1}^n \sum_{j=i+1}^n \frac{1}{m} \text{ (Explained above)} \\ &= \frac{\binom{n}{2}}{m} = \frac{n(n-1)}{2m} \text{ (There are } \binom{n}{2} \text{ total unordered pairs)} \end{aligned}$$

Hence the expected number of collisions is $\frac{n(n-1)}{2m}$.

Tutorial 6

1. Give an implementation for the scheduling problem using Heap.

Solution: Refer the slides to know about the scheduling problem.

See Algorithm 8 for the Pseudo-code.

Use of the Heap (Min Heap here) is to select the minimum remaining processing time among the processes present at any time. Before the timestamp of any new process is added, we will keep performing the jobs (available at that moment) with the minimum remaining times. As soon as a new process is added, the current process (if any) is preempted (i.e. interrupted) and the algorithm now looks for the process with the minimum remaining time. In the end, the for loop exits only when all the jobs are completed.

Given solution will have a time complexity of $\mathcal{O}(n \log n)$. Why? Think. (Hint: Number of operations (push or pop) won't exceed $4n$.)

Algorithm 8: Shortest Remaining Time Scheduling

Data: Array (sorted) *time* containing the timestamp the i^{th} job is called and a corresponding array *duration* containing the duration of the i^{th} job. Both are of size n .

```
1 Function Schedule( $n, time, duration$ ):
2   MinHeap RemTime // A Min Heap for the remaining processing times
3   for  $i \leftarrow 0$  to  $n - 1$  do
4     RemTime.push((duration[i], i)) // Sorted by the first of the pair
5     CurrTime  $\leftarrow time[i]$ 
6     if  $i = n - 1$  then
7       | EndTime  $\leftarrow \infty$  // Some large value
8     end
9     else
10      | EndTime  $\leftarrow time[i + 1]$ 
11    end
12    while not RemTime.empty() do
13      MinRemTime, CurrIndex  $\leftarrow$  RemTime.top()
14      RemTime.pop()
15      NextCurrTime  $\leftarrow$  CurrTime + MinRemTime
16      if NextCurrTime > EndTime then
17        | print("Job", CurrIndex, "starts at", CurrTime, "and ends at", EndTime)
18        | duration[CurrIndex]  $\leftarrow$  duration[CurrIndex] - (EndTime - CurrTime)
19        | RemTime.push((duration[CurrIndex], CurrIndex))
20        | break
21      end
22      print("Job", CurrIndex, "starts at", CurrTime, "and ends at", NextCurrTime)
23      print("Job", CurrIndex, "will be completed.")
24      CurrTime  $\leftarrow$  NextCurrTime
25    end
26  end
27 end
```

2. Can a Priority Queue be implemented as a red-black tree? What advantages does a heap implementation have over a red-black tree implementation?

Solution: Yes, the priority queue can be implemented using red-black trees. Priority queue has 3 operations which RB trees can also do:

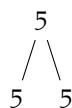
- insert: Insert can be done in an RB tree using the priority as the property of the tree
- top: Return the pointer to the rightmost element of the RB tree
- deleteMax: delete the rightmost element of the RB tree

Time complexity can be made equal for all of the above operations. (top can work in $\mathcal{O}(1)$ time. Think.)

The advantage of heap implementation is you can store the priority queue in an array, in which you aren't storing anything extra like child pointer, parent pointer etc, you are just storing the data of the node. Also, array implementation increases the cache performance.

3. Give a tree, if exists, that is a binary search tree, is a heap, and has more than two nodes. If such a tree does not exist, give a reason.

Solution: Consider the given binary tree:-



Since all keys are equal (and last level is complete hence no worry about being left-filled), this tree satisfies the definitions of both BST and Heap (Refer the slides for the definitions).

4. Suppose we have a 2D array where we maintain the following conditions: for every (i, j) , we have $A(i, j) \leq A(i + 1, j)$ and $A(i, j) \leq A(i, j + 1)$. Can this be used to implement a priority queue?

Solution: Yes, it is very similar to that of a heap. There are some modifications to the heap implementation:

- Store the negative of the priority in this data structure
- $A[i, j]$ has two children $A[i, j + 1]$ and $A[i + 1, j]$
- $A[i, j]$ has two parent $A[i, j - 1]$ and $A[i - 1, j]$
- Level of $A[i, j]$ is denoted by $i + j$
- You can again define notation of level-wise left filled.
- Deletion is done similar to heap
- Insertion is done similar to heap with a modification that the max of both parent is considered for swap.
- $A[0, 0]$ is the max-priority element.

5. *Given an unsorted array, find the k th smallest element using a priority queue.*

Solution: Here are the steps:-

1. Push all the elements of the array in the priority queue (Min Heap one).
2. Pop k times. Answer is the k^{th} pop.

Follow up: Can you do this in $\mathcal{O}(n)$ time (Not necessarily using Heaps)?

6. *Given two heaps give an efficient algorithm to merge the heaps.*

Solution: Here are the steps:-

1. Concatenate the array of the heaps (or copy all the elements from both arrays to a new one).
2. Apply BuildHeap function on this new array.

Time Complexity of this algorithm is $\mathcal{O}(m + n)$, m and n being the size of the heaps. This is because both the steps above are linear in time.

Tutorial 7

1. Compute KMP-table array h for pattern "babbaabba".

Solution: The algorithm at the end is given. In the beginning, index i is 1, and index j is 0. Note that the length m is 9 and the array h looks like this:

-1									
----	--	--	--	--	--	--	--	--	--

Now, $pat[1]$ is 'a' while $pat[0]$ is 'b'. Since the two are not equal, we assign $h[1]$ to 0 (j).

-1	0								
----	---	--	--	--	--	--	--	--	--

The while loop sets $j = h[j] = -1$. Then, both get incremented so that $i = 2$ and $j = 0$. Now, $pat[2] = pat[0]$ so we assign $h[2] = h[0] = -1$.

-1	0	-1							
----	---	----	--	--	--	--	--	--	--

Now, after incrementing and in the next iteration, $pat[3] \neq pat[1]$ so $h[3]$ is assigned to 1.

-1	0	-1	1						
----	---	----	---	--	--	--	--	--	--

One iteration of the (inner) while loop assigns j to $h[j]$ which is 0. Then, $pat[3] = pat[0]$ so it stops there. After incrementing both variables, we see that $pat[4] = pat[1]$ so we assign $h[4]$ to $h[1]$ which is 0.

-1	0	-1	1	0					
----	---	----	---	---	--	--	--	--	--

On incrementing both variables, $pat[5] \neq pat[2]$ so we assign $h[5]$ to 2.

-1	0	-1	1	0	2				
----	---	----	---	---	---	--	--	--	--

One iteration of the while loop assigns j to $h[j]$ which is -1. After incrementing, we get $i = 6$ and $j = 0$. Now, $pat[6] = pat[0]$ so $h[6]$ is assigned to $h[0]$ which is -1.

-1	0	-1	1	0	2	-1			
----	---	----	---	---	---	----	--	--	--

Next iteration, $pat[7] \neq pat[1]$ so $h[7]$ is assigned to 1.

-1	0	-1	1	0	2	-1	1		
----	---	----	---	---	---	----	---	--	--

The while loop resets j to -1, following which j is incremented to 0 and i to 8. Now, $pat[8] \neq pat[0]$ so $h[8]$ gets assigned to zero.

-1	0	-1	1	0	2	-1	1	0	
----	---	----	---	---	---	----	---	---	--

At last, $j = 0$ and $i = 9$. Now, since $i = m$ we break the standard loop and assign $h[i]$ to j , which is zero. At the end, the KMP table looks like this:

-1	0	-1	1	0	2	-1	1	0	0
----	---	----	---	---	---	----	---	---	---

Algorithm 9: KMPtable

```
1 Function KMPtable(String pat):
2    $i \leftarrow 1$ ;
3    $j \leftarrow 0$ ;
4    $m \leftarrow \text{pat.length}()$ ;
5    $h \leftarrow \text{new int } [n + 1]$ ;
6    $h[0] \leftarrow -1$ ;
7   while  $i < m$  do
8     if  $\text{pat}[i] \neq \text{pat}[j]$  then
9        $h[i] \leftarrow j$ ;
10      while  $j \geq 0$  and  $\text{pat}[i] \neq \text{pat}[j]$  do
11         $j \leftarrow h[j]$ ;
12      end
13    end
14    else
15       $h[i] \leftarrow h[j]$ ;
16    end
17     $i \leftarrow i + 1$ ;
18     $j \leftarrow j + 1$ ;
19  end
20   $h[m] \leftarrow j$ ;
21  return  $h$ ;
22 end
```

2. Is this version of the KMPtable algorithm correct?

Algorithm 10: KMPtableDiff

```
1 Function KMPtableDiff(String pat):
2    $i \leftarrow 1$ ;
3    $j \leftarrow 0$ ;
4    $m \leftarrow \text{pat.length}()$ ;
5    $h \leftarrow \text{new int } [n + 1]$ ;
6    $h[0] \leftarrow -1$ ;
7   while  $i < m$  do
8      $h[i] \leftarrow j$ ;
9     while  $j \geq 0$  and  $\text{pat}[i] \neq \text{pat}[j]$  do
10       $j \leftarrow h[j]$ ;
11    end
12     $i \leftarrow i + 1$ ;
13     $j \leftarrow j + 1$ ;
14  end
15   $h[m] \leftarrow j$ ;
16  return  $h$ ;
17 end
```

Solution: The answer is **no**. This algorithm is incorrect. Whenever $\text{pat}[i] = \text{pat}[j]$, the wrong value is inserted into the KMP table (meaning the table is incorrect). This was discussed in slides that if $h[i]$ is assigned to j , it does not meet the requirement that when the shifted string pat is aligned with pat , the compared indices i and j differ. If they do not, then the main search algorithm gets more inefficient.

3. Suppose that there is a letter z in pat of length n such that it occurs in only one place, say k , which is given in advance. Can we optimize the computation of h ?

Solution: Let the character z appear in position $m < n$ and nowhere else. Recall the function of the table is to find the largest index j given i such that $pat[0:j] = pat[i-j:i]$ but $pat[j] \neq pat[i]$. By definition, that j has to be less than i , else it implies a shift of the search string by 0 in the main search (since in case of inequality between $txt[i]$ and $pat[j]$, there, we move j to $h[j]$ and large $h[j]$ means small shift forward; $j = h[j]$ means no shift).

If $j > m$, it means that $pat[0:j]$ contains z . For the equality to hold, $pat[i-j:i]$ should contain z at index m , meaning that $pat[j+m] = z$, which goes against the premise of our question: there is only one index where z is located. So, we say that j in the computation can never exceed m or, more specifically, no entry in the KMP table can exceed m . Further, $pat[i-j:i]$ cannot contain z or index m in its range.

The computation of the KMP table till index m is the same as the earlier algorithm, and $h[m]$ will get assigned to the correct value anyway. For higher values, we split pat into $subpat = pat[0:m]$ and $src = pat[m+1:n]$. In order to compute the value of $h[i]$ for $i > m$, we try to follow the KMP search algorithm with src as the text and $subpat$ as its pattern. Except, we do not look for exact matches. Rather, when $src[i] \neq subpat[j]$, we assign $h[i+m+1] = j$, as expected.

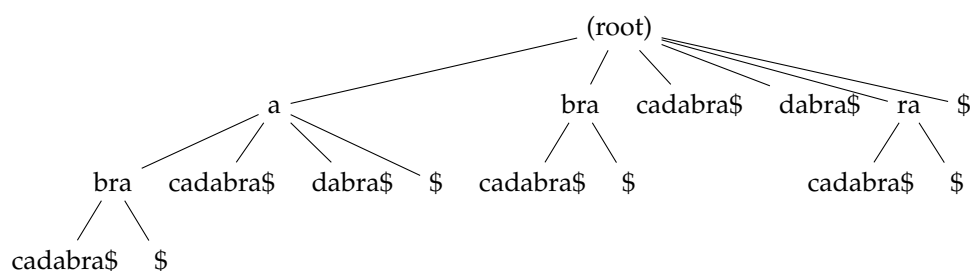
The case when $src[i] = subpat[j]$ seems a bit complicated, but just like the KMP table computation, doing $h[i+m+1] \leftarrow h[j]$ will work, since $src[i] = subpat[h[j]] \neq subpat[j]$ and all previous indices are equal. When $src[i] \neq subpat[j]$, after assigning $h[i+m+1] \leftarrow j$, we keep moving j much like KMPsearch and KMPtable algorithms.

This begs the question: is it really more efficient? To begin with, we cannot achieve less than $O(n)$ since we need to parse the entire string at least once (the output depends on all characters) and generate an $O(n)$ output. But can we make it faster $O(n)$? Such as reducing the constant c rather than focusing on asymptotic complexity?

One idea is to multi-thread the two computations, the KMP-table of $subpat$ and the modified KMP-search of $subpat$ in src that just updates the same table (algorithm not shown, since it gets cluttered and tedious with managing threads, and guaranteeing that i in the $subpat$ thread is not less than j in the KMPsearch – this requires locks and semaphores knowledge, and can still be as ‘inefficient’ which is $O(n)$). Any single-threaded solution (such as doing the 2 procedures in sequence or in one master function with somewhat on-demand computation of $h[j]$ for j in $subpat$) necessarily takes as much time as the normal KMPtable-filling algorithm.

4. Compute the suffix tree for ‘*abracadabra\$*’. Compress degree 1 nodes. Use substrings as edge labels. Put a square around nodes where a word ends. Use it to locate the occurrences of ‘*abr*’

Solution: This is the required suffix tree (dollar (\$) – where a word ends):



To locate 'abr', we traverse the suffix tree to node for 'a'. Then we look at the child corresponding to 'br'. It is the left child. We see that there are two suffixes: 'abra' (just that, ending in \$ symbol node) and 'abracadabra', which is the whole word itself.

5. Review the argument that for a given text T , consisting of k words, the ordinary trie occupies space, which is a constant multiple of $\|T\|$. How is it that the suffix tree for a text T is of size $O(\|T\|^2)$? Give a worst-case example.

Solution: Yes, for a given text, the trie occupies space a constant multiple of $n = \|T\|$, since each character has almost one node associated with it, and the space taken by each node in the trie is constant. Even if the alphabet is infinite, a different representation of trie (linked list for children) implies that each trie-node has at max one parent, and therefore, the space taken is directly proportional to the number of nodes.

On the other hand, the suffix tree stores **all** suffixes. Consider the text T having only one word: each suffix (corresponding to a character in the word, of the n) takes up $O(\text{size of suffix})$ space. The net space evaluates to $O(n^2)$ since suffixes are substrings of length $1, 2, \dots, n$ where n is the length of one word.

Sure, there may be patterns that are repeated (and therefore, the number of nodes is fewer), and the best case is $O(n)$ for strings like $aaaaaa \dots aa$. However, the worst case is $O(n^2)$ when all characters are distinct (particularly when there is no restriction on character space), and the text comprises one word.

Tutorial 8

1. In an Huffman code instance, show that if there is a character with frequency greater than $2/5$, then there is a codeword of length 1. Further, show that if all frequencies are less than $1/3$ then there is no codeword of length 1.

Solution: According to Huffman compression, the two trees that have the least net frequency become the left and right child of a new node, forming a new tree. This process is repeated till there is only one tree left. Recall that in Huffman compression, a tree can have one node (root and leaf) that represents one character, or it can be any binary tree, such that each leaf represents a character and the net frequency of the tree is the sum of frequencies of occurrences of all the leaves (characters) in the text.

1. There exists a frequency more than $2/5$, then there exists a codeword of length 1.

If that character c_0 is not the most frequent character, but it has $f_{c_0} > 2/5$, then there exist two or more characters with frequencies $2/5$ or higher. Let the second be c_1 with $f_{c_1} \geq f_{c_0} > 2/5$. If this is the case, either there are just c_0 and c_1 (they are encoded 0 or 1), or there are more than these two characters, but all the rest have net frequency $1 - f_{c_0} - f_{c_1} < 1/5$. All of those get clubbed together to form one tree, since, until this is the case, the two trees with the least frequency cannot include either of c_0 and c_1 .

After this happens, the 'less frequent' tree T gets clubbed with c_0 first, then the resultant gets clubbed with c_1 [Clubbed = they become children of a new node, forming one tree]. At the end, the c_1 is either the left or right child of the root of the final Huffman tree, and has a codeword of length 1 (codewords '0' or '1')

If c_0 is the most frequent character and we expect it to not have a codeword of length 1, it should be clubbed with another tree sometime before the last step, when there are atleast two trees, both of which do not contain c_0 . Further, all trees except for $\{c_0\}$ and one other tree (that gets clubbed with c_0 whenever it gets clubbed), should have a net frequency greater than or equal to $f_{c_0} > 2/5$. This is possible only when there is one such tree, and there are three trees in all.

Let us examine that case: let T_0 be the tree with highest frequency and T_1 the low frequency tree, that gets clubbed with c_0 . Clearly, $f_{T_0} > 2/5$ from this and the frequency of c_0 , we say that $f_{T_1} < 1/5$. Since T_0 is not a singleton tree, we look at the previous step wherever T_0 was formed from two children, T_2 and T_3 . Without loss of generality, let $f_{T_2} \geq f_{T_3}$.

Then, $2 \cdot f_{T_2} \geq f_{T_2} + f_{T_3} = f_{T_0} > 2/5$ so, before its combination, $f_{T_2} > 1/5 > f_{T_1}$; there existed atleast two trees T_3 and T_1 (more if T_1 was not formed yet), both with net frequency less than that of T_2 . This means that if T_2 was clubbed (it had to be for our case to hold), its combination with T_2 did not follow the Huffman algorithm. This arises a contradiction.

Therefore, if there exists a character of frequency more than $2/5$, then there has to be a character with codeword of length 1.

2. If all frequencies are less than $1/3$, there exists no codeword of length 1.

Let there be a character c_0 that has the luxury of a codeword of length 1. If this were so, at the last iteration of Huffman, we have a singleton tree c_0 (else, it would not have a codeword of length 1), and a tree T_0 . Since $f_{c_0} < 1/3$, $f_{T_0} > 2/3$.

Since this tree T_0 has to contain two or more characters, it must have two subtrees T_2 and T_3 , and those must have been clubbed together in the last step (the only possibility). Without loss of generality, let $f_{T_2} \geq f_{T_3}$. This implies that $2 \cdot f_{T_2} \geq f_{T_2} + f_{T_3} = f_{T_0} > 2/3$

or $f_{T_2} > 1/3 > f_{c_0}$. Then, the two trees with least net frequencies are T_3 and $\{c_0\}$. Those would have been clubbed together, arising in a contradiction.

Therefore, we conclude that there cannot exist a codeword of length 1 if all characters have frequency less than $1/3$.

2. Suppose that there is a source that has three characters 'a', 'b', 'c'. The output of the source cycles in the order of 'a', 'b', 'c' followed by 'a' again, and so on. In other words, if the last output was a 'b', then the next output will either be a 'b' or a 'c'. Each letter is equally probable. Is the Huffman code the best possible encoding? Are there any other possibilities? What would be the pros and cons of this?

Solution: Yes, we can use a location-aware encoding i.e. one wherein the encoding of a character depends on its surroundings, specifically, the character before it. This is a better encoding scheme than Huffman encoding.

For instance, we choose to encode a character with a '1', assuming that it is different from the previous character (there can only be one possible value corresponding to this 1), and a '0' assuming it is same. We can encode the first character using any Huffman tree, say, 1 for 'a', 00 for 'b' and 01 for 'c'. Then, abccaaabbca will be encoded as 11101001011. Also, 010011010110 would correspond to cccabbccabb.

Roughly, this corresponds to one bit per character. Compare that with using on an average $5/3$ bits if we encoded each character individually (and they had roughly same frequency each; true since each letter is equally probable), which would be the case with Huffman encoding. This clearly means that our encoding is better, the big pro. Also, ours does not take a very long time to evaluate (definitely not more than Huffman). Another pro (plus point) is that this is very much scalable to m character systems, say the English alphabet, except, when we know that an 'a' is followed by 'a' or 'b'; 'b' by 'b' or 'c'; ... and 'z' by 'z' or 'a'.

Speaking of cons (downsides) of our encoding, this is not generalizable to arbitrary strings. If 'ac', 'cb' or 'ba' occurred as a substring anywhere (very likely), then our encoding fails (there does not exist an encoding by our scheme). Granted, we know that this is not the case for our source; but we first need to verify that the source is of this form (this can easily be done) so this is not really a con.

Another improvement would be to use something like run-length-encoding algorithm, particularly efficient for strings with long runs like aaaaaaaaaabbbbbbbbbbbbbbcccccccc; but those would be inefficient if there are lot of short runs, particularly, if lot of runs have length 1.

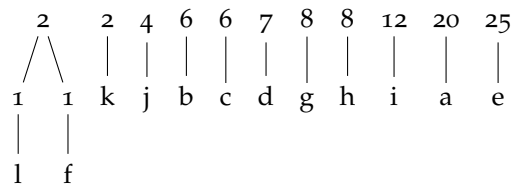
3. Given the following frequencies, compute the Huffman code tree.

Character	a	d	g	j	b	e	h	k	c	f	i	l
Frequency	20	7	8	4	6	25	8	2	6	1	12	1

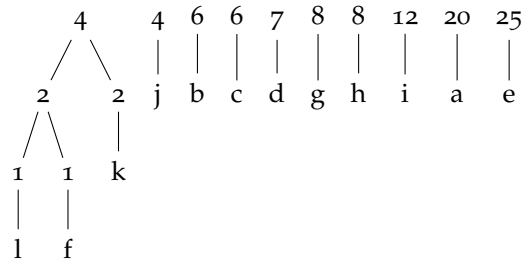
Solution: The initial forest for Huffman algorithm is:

1	1	2	4	6	6	7	8	8	12	20	25
l	f	k	j	b	c	d	g	h	i	a	e

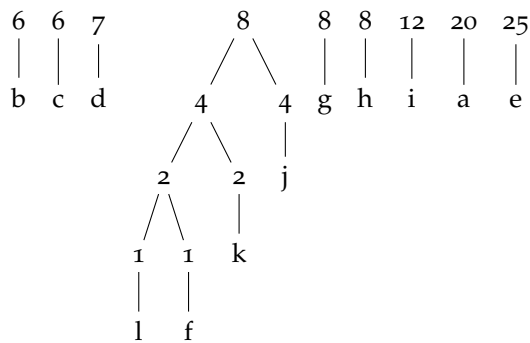
The first iteration combines the two trees with least frequency (l and f), resulting in our forest looking like this:



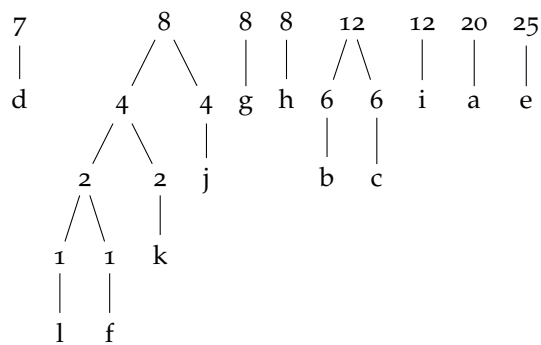
After the next iteration, our forest looks like:



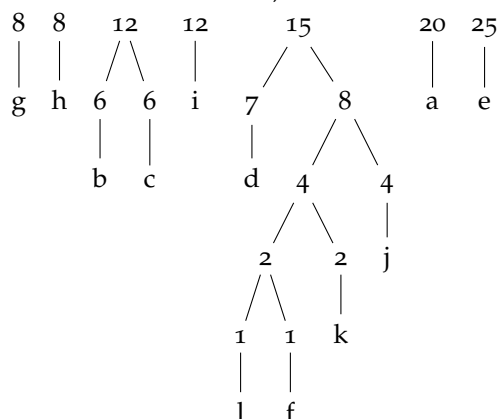
After the next iteration, our forest looks like:



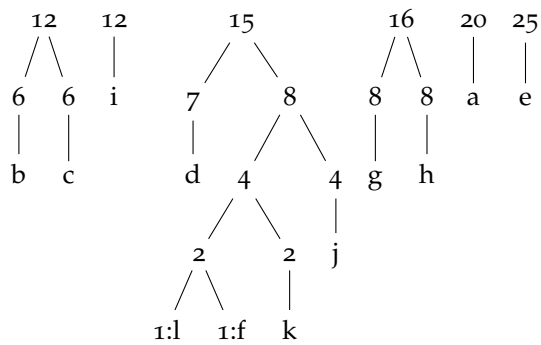
After the next iteration, our forest looks like:



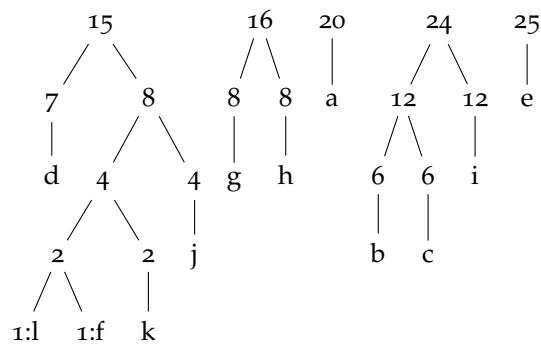
After the next iteration, our forest looks like:



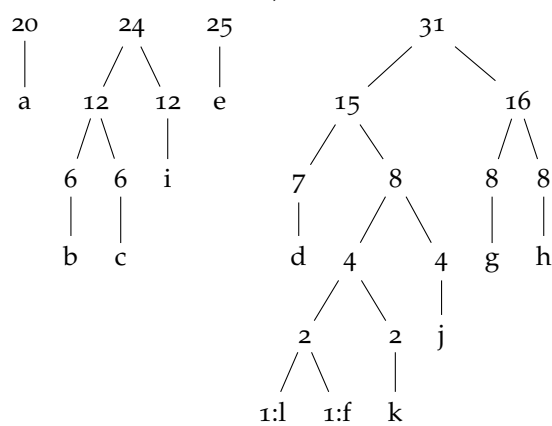
After the next iteration, our forest looks like:



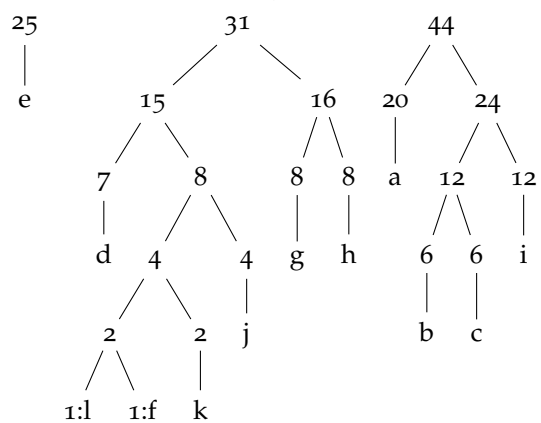
After the next iteration, our forest looks like:



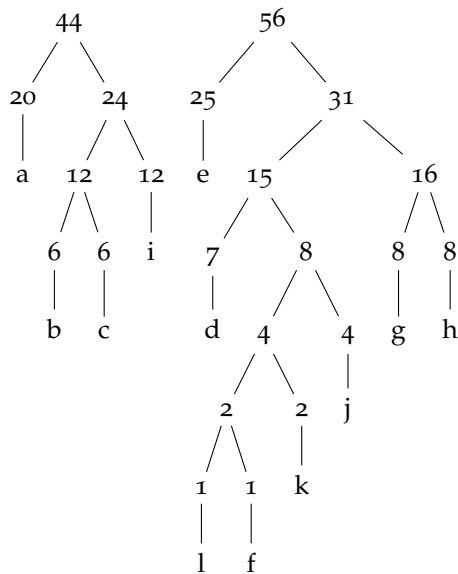
After the next iteration, our forest looks like:



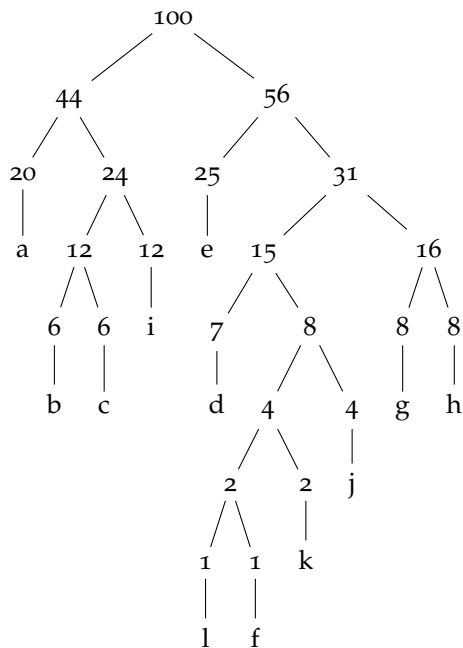
After the next iteration, our forest looks like:



The penultimate iteration involves:



... at the end, we have our complete Huffman tree!



4. We have seen a “proof of optimality” of the Huffman Coding. Is Huffman Coding really the best lossless text compression method out there? If no, give an example where Huffman Coding costs more than some other coding. Then is Huffman “optimal”? What “optimality” did we talk back there? Can you think why is Huffman not the best then?

Solution: No, Huffman Coding is NOT the best text compression method we have.

There will be 2 cases:- 1) We build a different Huffman tree for every string, 2) We use some universal Huffman tree for all strings.

In case 1, we need to use some bits to store the Huffman tree made every time for an input string. Take “Today_is_Monday” for example. Original string takes $15 \times 8 = 120$

bits (Assuming 8 bits for each ASCII), but Huffman takes *atleast* $49 + 10 \times 8 = 129$ bits (49 for the Huffman string created (You can work this out on your own) and the table (which maps each character of the word to a Huffman code) will contain 10 entries as 10 distinct characters (Also assumed that each entry takes minimum 8 bits)).

In case 2, if there is some pre-processed Huffman table, then take “xxxx”. Since ‘x’ is the least used English alphabet, it’s Huffman code will likely take more than 8 bits.

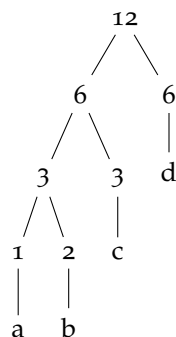
Huffman is “optimal” only when we consider bits used to represent the Huffman string. The proof done in class doesn’t cover the bits that maybe required for making the Huffman table.

Another reason why Huffman isn’t the best is that it assigns *integer*-number of bits to any symbol. Huffman codes are constrained to bit boundaries. In a Shannon information theoretic sense, if we just look at context-free probabilities of symbols, the ideal bit-length of a symbol is $-\log_2(P(x))$ where $P(x)$ is the probability of occurrence of the symbol x .

That number is unlikely to be an integer! The Huffman code’s inefficiency lies in the distance between the information theoretic optimal bit length and the actual integer bit length the Huffman code assigned.

5. What is the largest Huffman code for any character that can be created in a word? Give an example.

Solution: Largest Huffman code for any character can be one less than number of distinct characters in the word. Consider “abbccddddd”:-



See that code length for ‘a’ or ‘b’ is 3.

6. You are given a stick of length l . Your job is to break the stick to create pieces of lengths $l_1, l_2, l_3, \dots, l_m$ (These sum up to l). But the cost of breaking any stick (at any point) into two parts is equal to the length of that stick. What will be the minimum cost to do your job?

Solution: The solution is pretty similar to that of Huffman Coding. We build a tree just like the Huffman tree, just here the weights of the nodes is the length of the sticks instead of frequency of character. Proof of optimality follows from that of Huffman’s.