



KIET Group of Institutions, Ghaziabad

Department of Computer Applications

(An ISO – 9001: 2015 Certified & 'A' Grade accredited Institution by NAAC)



NOTES

Subject Name: Data Structures & Analysis of Algorithms **Subject Code:** KCA-205 **Semester:** II

UNIT-4

Basic terminology used with Tree, Binary Trees, Binary Tree Representation: Array Representation and Pointer (Linked List) Representation, Binary Search Tree, Complete Binary Tree, A Extended Binary Trees, Tree Traversal algorithms: Inorder, Preorder and Postorder, Constructing Binary Tree from given Tree Traversal, Operation of Insertion, Deletion, Searching & Modification of data in Binary Search Tree. Threaded Binary trees, Huffman coding using Binary Tree, AVL Tree and B Tree.

3.1 Basic Terminologies

Terminologies used in Trees

- **Root** – The top node in a tree.
- **Child** – A node directly connected to another node when moving away from the Root.
- **Parent** – The converse notion of a *child*.
- **Siblings** – Nodes with the same parent.
- **Descendant** – A node reachable by repeated proceeding from parent to child.
- **Ancestor** – A node reachable by repeated proceeding from child to parent.
- **Leaf** – A node with no children.
- **Internal node** – A node with at least one child.
- **External node** – A node with no children.
- **Degree** – Number of sub trees of a node.
- **Edge** – Connection between one node to another.
- **Path** – A sequence of nodes and edges connecting a node with a descendant.
- **Level** – The level of a node is defined by $1 +$ (the number of connections between the node and the root).
- **Height of node** – The height of a node is the number of edges on the longest downward path between that node and a leaf.

- **Height of tree** – The height of a tree is the height of its root node.
- **Depth** – The depth of a node is the number of edges from the node to the tree's root node.
- **Forest** – A forest is a set of $n \geq 0$ disjoint trees.

3.2 Definition and Representation

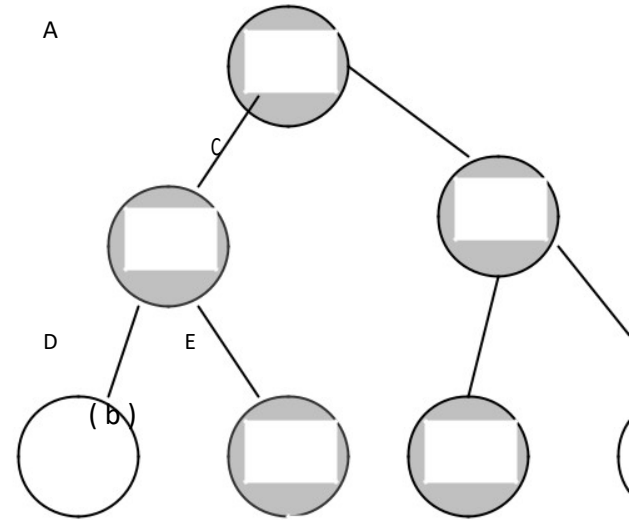
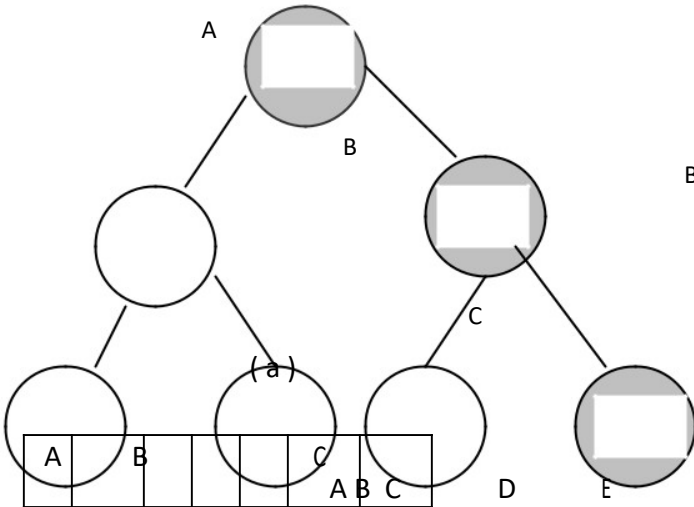
A tree is a (possibly non- linear) data structure made up of nodes or vertices and edges without having any cycle. The tree with no nodes is called the **null** or **empty** tree. A tree that is not empty consists of a root node and potentially many levels of additional nodes that form a hierarchy.

3.3 Representation of Binary Tree

There are two representations used to implement binary trees.

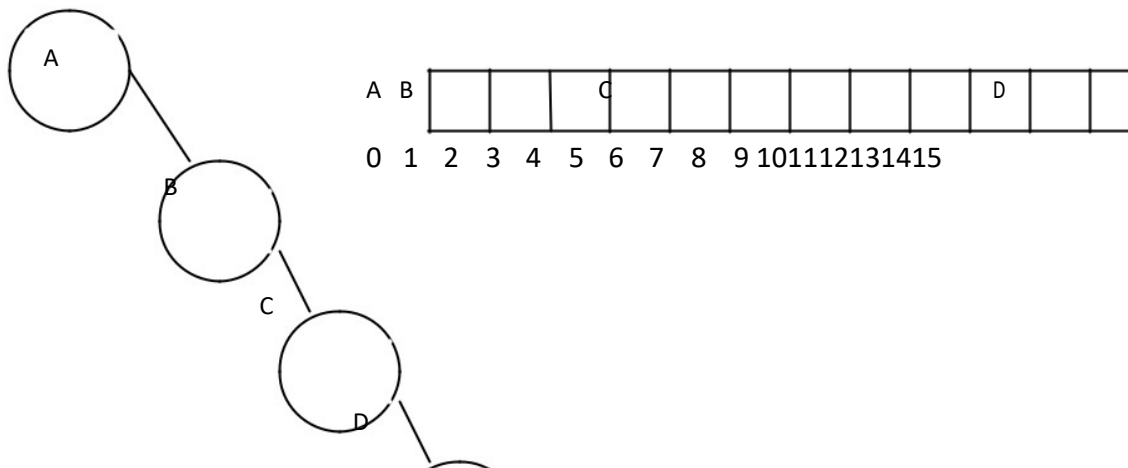
(i) Array Representation (ii) Linked list Representation

Array Representation: In this, the given binary trees even though are not complete binary trees, they are shown as complete binary trees in which missing elements are un shaded circles. The array representations for the following trees are shown in below.



In array, the elements of the binary are placed in the array according to their number assigned. The array starts indexing from 1. The main drawback of array representation is wasteful of memory when there are many missing elements.

The binary tree with n elements requires array size up to 2^n . Suppose array positions indexing from 0, then array size reduces to $2^n - 1$. The right skewed binary trees have maximum waste of space. The following right skewed binary tree's array representation is shown as follows.



Linked list Representation: The most popular way to represent a binary tree is by using links or pointers. The node structure used in this representation consists of two pointers and an element for each node. The node structure is given as:

leftchild	element	rightchild
-----------	---------	------------

The first field leftchild pointer maintains left child of it. The middle field is the element of the node and last field is the right child pointer maintains right child of it.

Binary tree traversals: There are four ways to traverse a binary tree. They are

(a) Pre order (b) In order (c) Post order (d) Level order

The first three traversals are performed using recursive approach and are done using linked list scheme. In these, the left sub tree is visited before visiting right sub tree. The difference among these is position of visiting the node.

(a) **Pre order:**

i) Visit Root node. ii) Visit Left sub tree. iii) Visit Right Sub tree.

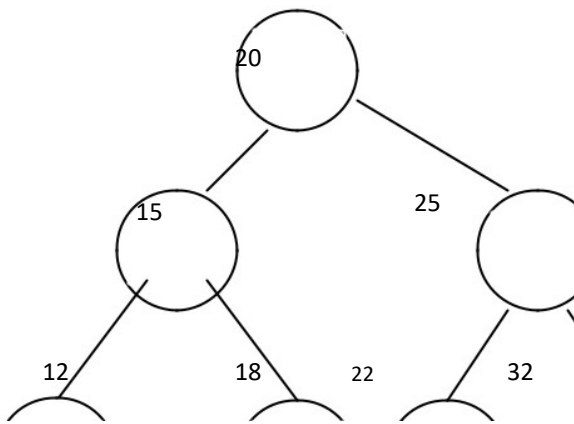
(b) **In order:**

i) Visit Left sub tree. ii) Visit Root node. iii) Visit Right Sub tree.

(c) **Post order:**

i) Visit Left sub tree. ii) Visit Root node. iii) Visit Right Sub tree

The following is an example binary tree with pre order, in order, post order and level order traversals:



(a)

pre order is : 20 15 12 18 25 22 32

in order is : 12 15 18 20 22 25 32

post order is : 12 18 15 22 32 25 20

level order is : 20 15 25 12 18 22 32

Tree Traversal

Traversal is a process to visit all the nodes of a tree and may print their values too. Because, all nodes are connected via edges (links) we always start from the root (head) node. That is, we cannot randomly access a node in a tree. There are three ways which we use to traverse a tree –

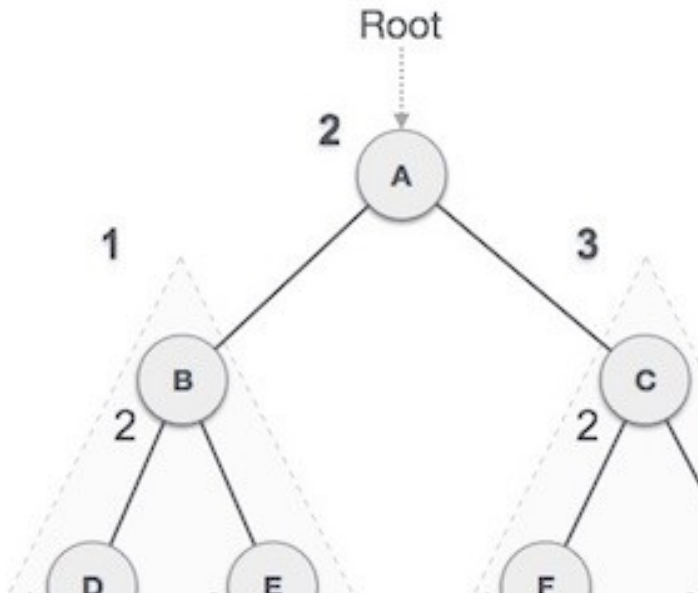
- In-order Traversal
- Pre-order Traversal
- Post-order Traversal

Generally, we traverse a tree to search or locate a given item or key in the tree or to print all the values it contains.

In-order Traversal

In this traversal method, the left subtree is visited first, then the root and later the right sub-tree. We should always remember that every node may represent a subtree itself.

If a binary tree is traversed **in-order**, the output will produce sorted key values in an ascending order.



We start from **A**, and following in-order traversal, we move to its left subtree **B**. **B** is also traversed in-order. The process goes on until all the nodes are visited. The output of inorder traversal of this tree will be –

$D \rightarrow B \rightarrow E \rightarrow A \rightarrow F \rightarrow C \rightarrow G$

Algorithm

Until all nodes are traversed –

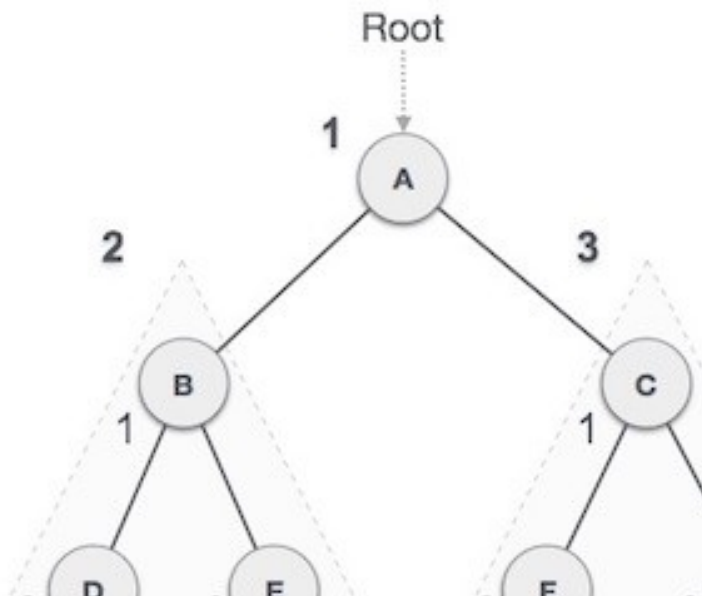
Step 1 – Recursively traverse left subtree.

Step 2 – Visit root node.

Step 3 – Recursively traverse right subtree.

Pre-order Traversal

In this traversal method, the root node is visited first, then the left subtree and finally the right subtree.



We start from **A**, and following pre-order traversal, we first visit **A** itself and then move to its left subtree **B**. **B** is also traversed pre-order. The process goes on until all the nodes are visited. The output of pre-order traversal of this tree will be –

$A \rightarrow B \rightarrow D \rightarrow E \rightarrow C \rightarrow F \rightarrow G$

Algorithm

Until all nodes are traversed –

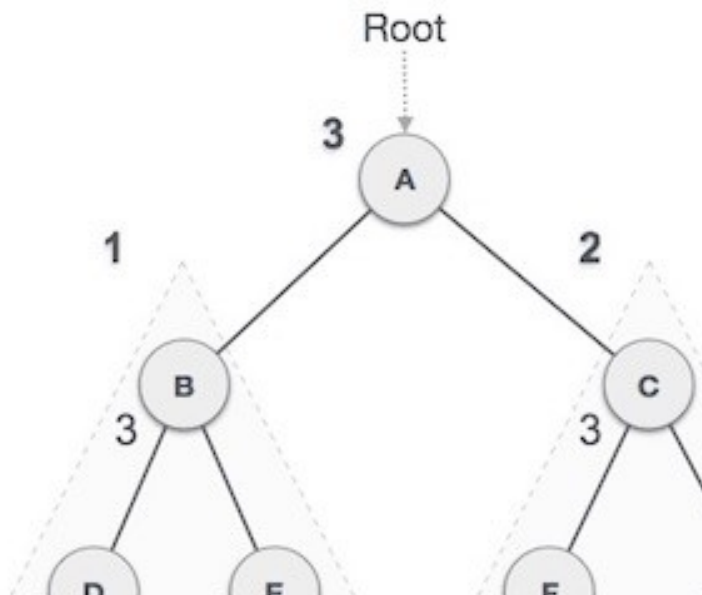
Step 1 – Visit root node.

Step 2 – Recursively traverse left subtree.

Step 3 – Recursively traverse right subtree.

Post-order Traversal

In this traversal method, the root node is visited last, hence the name. First we traverse the left subtree, then the right subtree and finally the root node.



We start from **A**, and following Post-order traversal, we first visit the left subtree **B**. **B** is also traversed post-order. The process goes on until all the nodes are visited. The output of post-order traversal of this tree will be –

$$D \rightarrow E \rightarrow B \rightarrow F \rightarrow G \rightarrow C \rightarrow A$$

Algorithm

Until all nodes are traversed –

Step 1 – Recursively traverse left subtree.

Step 2 – Recursively traverse right subtree.

Step 3 – Visit root node.

Tree Construction

We are given the Inorder and Preorder traversals of a binary tree. The goal is to construct a tree from given traversals.

Inorder traversal – In this type of tree traversal, a left subtree is visited first, followed by the node and right subtree in the end.

Inorder (tree root)

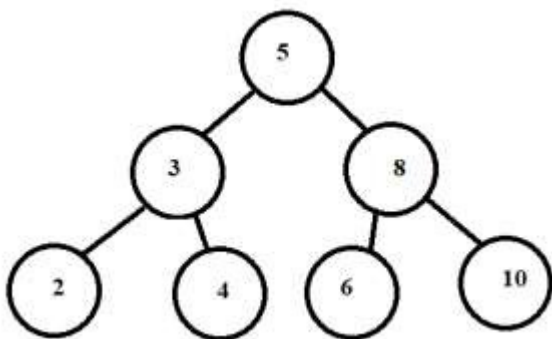
- Traverse left subtree of node pointed by root, call inorder (root→left)
- Visit the root
- Traverse right subtree of node pointed by root, call inorder (root→right)

Preorder traversal – In this type of tree traversal, the node visited first, followed by the left subtree and right subtree in the end.

Preorder (tree root)

- Visit the root
- Traverse left subtree of node pointed by root, call inorder (root→left)
- Traverse right subtree of node pointed by root, call inorder (root→right)

The inorder and preorder traversal of below tree are given –



Inorder

2-3-4-5-6-8-10

Preorder

4-3-2-5-8-6-10

Now we'll construct the above tree again for given preorder and inorder traversals.

Inorder

2-3-4-5-6-8-10

Preorder

5-3-2-4-8-6-10

- As we know that preorder visits the root node first then the first value always represents the root of the tree. From above sequence 5 is the root of the tree.

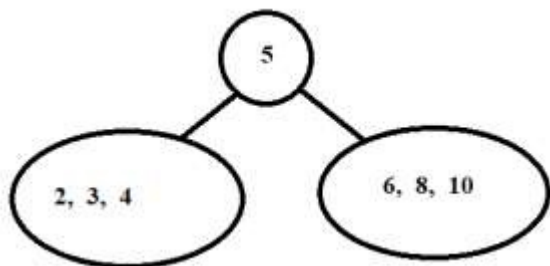
Preorder

5 -3-2-4-8-6-10

- From above inorder traversal, we know that a node's left subtree is traversed before it followed by its right subtree. Therefore, all values to the left of 5 in inorder belong to its left subtree and all values to the right belong to its right subtree.

Inorder

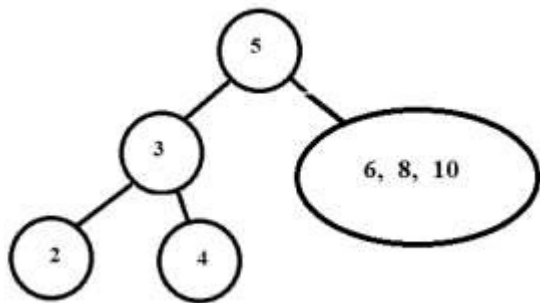
2-3-4 \leftarrow 5 \rightarrow 6-8-10



- Now for the left subtree do the same as above.

Preorder traversal of left subtree is 3 -2-4. So 3 becomes the root.

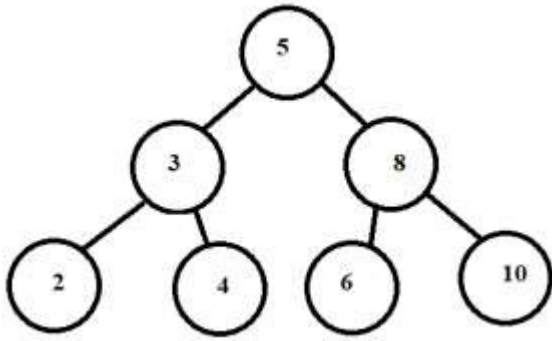
Inorder traversal divided further into $2 \leftarrow 3 \rightarrow 4$



- Now for the right subtree do the same as above.

Preorder traversal of the right subtree is 8 -6-10. So 8 becomes the root.

Inorder traversal divided further into $6 \leftarrow 8 \rightarrow 10$



So, in this way we constructed the original tree from given preorder and inorder traversals.

3.4 Types of Binary trees

3.4.1 Binary Search Trees:

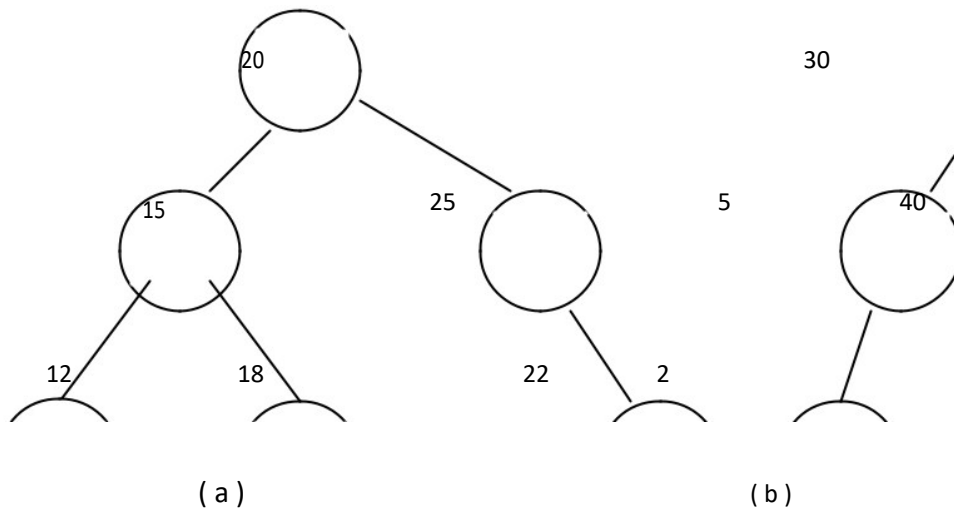
Any empty binary tree is an Binary Search Tree. A nonempty binary search tree has the following properties.

- (i) Every element has the key or value and no two elements have the same key. Therefore, all keys in the tree must be distinct.

- (ii) Any element key in left sub tree is less than the key of the root.
- (iii) Any element key in right sub tree is greater than the key of the root.
- (iv) Both left and right sub trees are also binary search trees.

There is some redundancy in this definition. The properties 2, 3, 4 together imply the keys must be distinct.

Some binary trees in which the elements with distinct keys are shown in the following figures.



The number inside a node is the element key. The tree (a) is not a binary search tree because the right sub tree of element 25 violating property 4. It means 22 is smaller than its parent 25. The trees of (b) is not a binary search tree.

When the property all keys are distinct is removed, then property 2 is replaced by smaller or equal and property 3 is replaced by larger or equal. The resulting tree is called a binary search tree with duplicates.

3.4.2 Binary Search tree Applications: It is mainly used in

- (i) Histogramming
- (ii) Best fit bin packaging
- iii) Crossing Distribution

3.5.2 Heap Trees

In computer science, a **heap** is a specialized tree-based data structure that satisfies the *heap property*: If A is a parent node of B then the key of node A is ordered with respect to the key of node B with the same ordering applying across the heap. A heap can be classified further as either a "**max heap**" or a "**min heap**". In a max heap, the keys of parent nodes are always greater than or equal to those of the children and the highest key is in the root node. In a min heap, the keys of parent nodes are less than or equal to those of the children and the lowest key is in the root node. Heaps are crucial in several

efficient graph algorithms such as Dijkstra's algorithm, and in the sorting algorithm heap sort. A common implementation of a heap is the binary heap, in which the tree is a complete binary tree (see figure).

In a heap, the highest (or lowest) priority element is always stored at the root, hence the name **heap**. A heap is not a sorted structure and can be regarded as partially ordered. As visible from the heap-diagram, there is no particular relationship among nodes on any given level, even among the siblings. When a heap is a complete binary tree, it has a smallest possible height—a heap with N nodes always has $\log N$ height. A heap is a useful data structure when you need to remove the object with the highest (or lowest) priority.

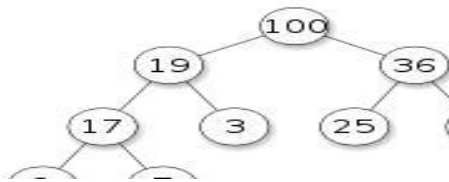


Fig (i) Heap Tree

Note that, as shown in the graphic, there is no implied ordering between siblings or cousins and no implied sequence for an in-order traversal (as there would be in, e.g., a binary search tree). The heap relation mentioned above applies only between nodes and their parents, grandparents, etc. The maximum number of children each node can have depends on the type of heap, but in many types it is at most two, which is known as a binary heap.

The heap is one maximally efficient implementation of an abstract data type called a priority queue, and in fact priority queues are often referred to as "heaps", regardless of how they may be implemented. Note that despite the similarity of the name "heap" to "stack" and "queue", the latter two are abstract data types, while a heap is a specific data structure, and "priority queue" is the proper term for the abstract data type.

A *heap* data structure should not be confused with *the heap* which is a common name for the pool of memory from which dynamically allocated memory is allocated. The term was originally used only for the data structure.

3.6 Height Balanced Trees

Height balanced trees (or AVL trees) is named after its two inventors, G.M. Adelson-Velskii and E.M. Landis, who published it in their 1962 paper "An algorithm for the organization of information." As the name suggests AVL trees are used for organizing information.

AVL trees are used for performing search operations on high dimension external data storage.

For example, a phone call list may generate a huge database which may be recorded only on external hard drives, hard-disks or other storage devices.

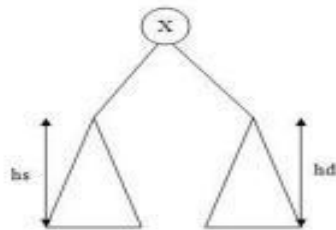
AVL is a data structure that allows storing data such that it may be used efficiently regarding the operations that may be performed and the resources that are needed. AVL trees are binary search trees, which have the balance propriety. The balance property is true for any node and it states: “the height of the left subtree of any node differs from the height of the right subtree by 1”.

The structure of the nodes of a balanced tree can be represented like:

```
struct NodeAVL{  
    int key;  
    int ech;  
    node *left, *right;  
};
```

Where: - key represents the tag of the node(integer number),

ech represents the balancing factor - left and right represent pointers to the left and right children.



Here are some important notions:

[1] The length of the longest road from the root node to one of the terminal nodes is what we call the height of a tree.

[2] The difference between the height of the right sub tree and the height of the left sub tree is what we call the balancing factor.

[3] The binary tree is balanced when all the balancing factors of all the nodes are -1,0,+1. We present below the function drum, which calculates the longest road from the current node , meaning the height of a sub tree:

```
void path(NodeAVL* p,int &max, int length)  
{  
    if (p!=NULL)  
    {  
        path(p->right,max,length+1);  
        if((p->left==NULL)&&(p->right==NULL)&&(maxleft,max,length+1);  
    }
```

```
}
```

Using this function we can determine the balancing indicator of each node of the tree with the function balance Factor:

```
void balanceFactor(NodeAVL *p)
{
    int max1,max2;
    max1=1;
    max2=1;
    if(p->left!=NULL)
        path(p->left,max1,1);
    else max1=0;
    if(p->right!=NULL)
        path(p->right,max2,1);
    else max2=0;
    p->ech=max2- max1;
}
```

3.7 B Trees

AVL and red black trees are used when the dictionary is small enough to reside in internal memory.

The search trees of higher degree are needed to get better performance for external dictionaries. ISAM is used to get good sequential and random access for external dictionaries.

3.8 m-way search trees: An m- way search tree may be empty. If it is not empty, it must satisfy the following properties.

- (i) In the tree, each internal node has up to m children and has elements between 1 & m-1.
- (ii) Every internal node with p elements has p+1 children.
- (iii) Consider any node with p elements. Let k_1, k_2, \dots, k_p be the keys of these elements. The elements are ordered such that $k_1 < k_2 < \dots < k_p$. Let c_0, c_1, \dots, c_p be p+1 children of the node. The elements in the sub tree with root c_0 have keys smaller than k_1 , those in the sub tree with root c_p have keys larger than k_p and those in sub tree with root c_i have the keys larger than k_i but smaller than k_{i+1} where $i \leq p$. Although external nodes are included when defining m-way search tree, external nodes are not represented physically in actual representation.

Consider the following seven way search tree, which have external nodes as solid squares and all other nodes are internal nodes. The root has 2 elements and 3 children. The middle child of the root has 6 elements and 7 children in which 6 are external nodes.

3.9.1.1 searching: To search for an element with key 31, begin checking with root first. The searching drops to middle child of the root because 31 lies between 10 & 80. Since search finds $k_2 (=30) < 31 < k_3 (=40)$, search drops to third sub tree of this node. There $31 < k_1 (=32)$, search moves to first sun tree of this node but external node is reached. When search terminates at external node, the key is not found. Otherwise means search exited at any internal node, the key is found.

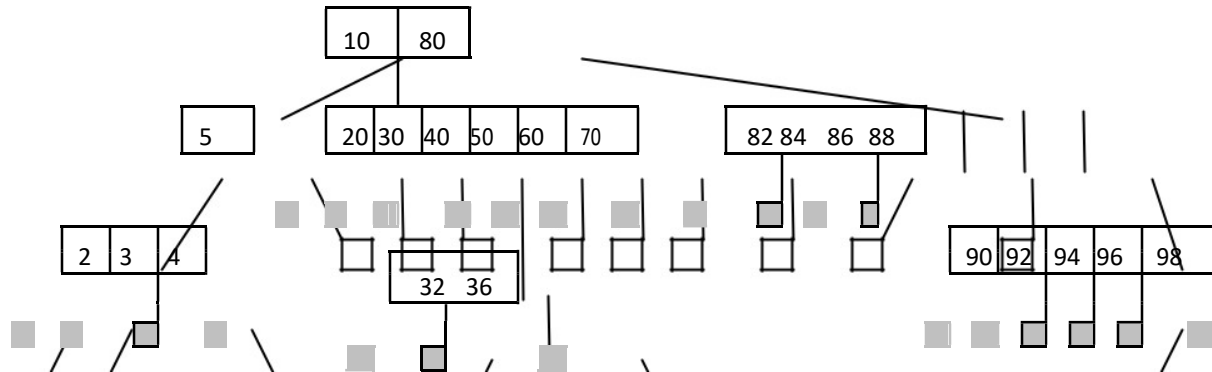


Fig:1 m-way search tree example

3.9.1.2 insertion: To insert an element with key, first search the tree. If it doesn't contain the key, then insert it. To insert the key 31, the search begins at root, then goes to middle child of root, then third child of this node, then first child of this node which is the external node. Since the node [32,36] can hold up to 6 elements, the new element is inserted as first element of this node.

Another example is insert an element with key 65, the search terminates at sixth child of middle child of the root. The new element is created and inserted into there.

3.9.1.3 deletion: To delete an element with key, first search for it. If it is there, then delete it. (i) If the deleted element is 20 in fig above, search for it. The searching ends at first element of the middle child of the root. Since its childrens c_0 & c_1 are 0, it can be deleted easily and results a node [30,40,50,60,70].

(ii) To delete 84, search ends at second element in third child of the root. Since its childrens c_1 & c_2 are 0, it can be deleted easily by resulting [82,86,88].

(iii) To delete an element with key 5, more work to be done. It has a nonnull c_0 and c_1 is external node. The largest key in the c_0 is brought to the deleted node place.

(iv) To delete an element with key 10, the root take either largest in its c_0 or smallest in c_1 . Suppose 5 of c_0 was brought to the root, 4 in the c_0 of deleted node 5 is brought to 5's old place.

3.9.1.4 Height: An m- way search tree of height h may have as few as h elements and as many as m^{h-1} .

This upper bound is obtained from the levels 1 through $h-1$ has exactly m children and nodes at level h have no children. Since each of these nodes has $m-1$ elements, the number of elements is m^{h-1} . An

200 way search tree of height 5 can hold $32 * 10^{10} - 1$ elements but might fold as few as 5 only.

3.9.2 B – trees: A B – tree of order m is an m – way search tree. If the B – tree is not empty, the corresponding extended tree must satisfy the following properties:

- (i) The root has at least two children.
 - (ii) All internal nodes other than the root have at least $\lceil m/2 \rceil$ children.
 - (iii) All external nodes are at the same level.
- (i) all the external nodes are not on the same level.
- (ii) Some of the internal nodes have two (= node [5]) and three (= node 32,36)) children which is not satisfying 2^{nd} property i.e $\lceil 7/2 \rceil = 4$ children.

The following is an B-tree of order 7:

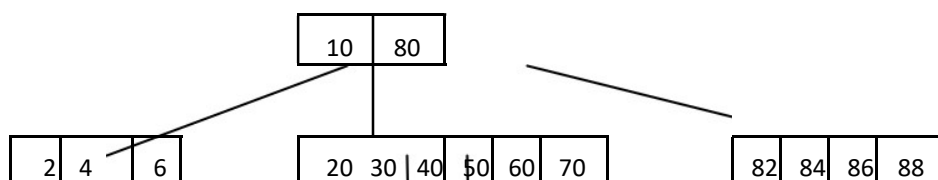


Fig 2 : B – tree of order 7

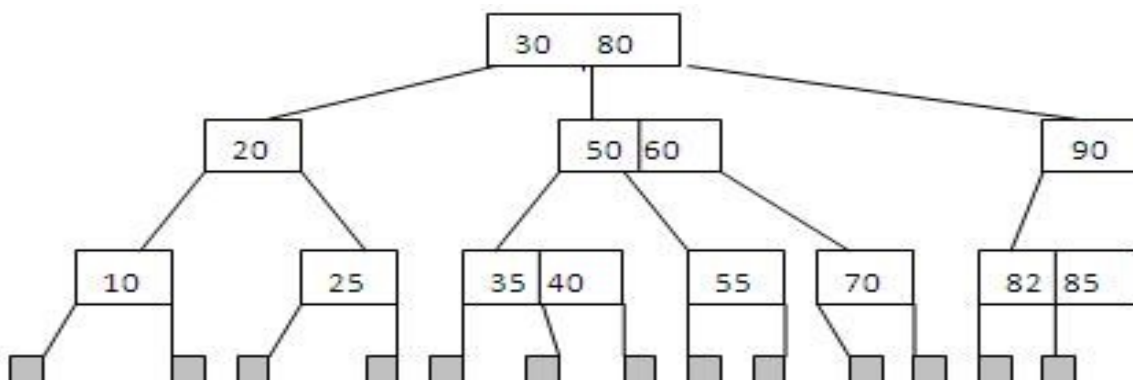
In a B-tree of order 2, all the internal nodes have exactly two children. This requirement coupled with all external nodes on the same level results full binary trees.

In a B-tree of order 3, internal nodes have either two or three children. It is also called 2-3 tree.

In a B- tree of order 4, internal nodes have two, three or four children. These are also referred as 2-3-4 trees and are called 2-4 trees. The following is an 2-3 tree. It becomes 2-3-4 tree when adding 14 and 16 to left child of 20.

3.9.2.1 Height of a B – tree: Let T be a B – tree of order m and height be h. Let $m = \lceil d/2 \rceil$ and n be number of elements in T.

- (a) $2d^{h-1} - 1 \leq n \leq m^h - 1$
- (b) $\log_m (n+1) \leq h \leq \log_d (n+1/2) + 1$



3.9.2.1 Searching: The searching an element in a B-tree is same as algorithm used for m-way search tree. Searching an element in an internal node of a B-tree of height takes at most h because all internal nodes need to be checked during the search.

3.9.2.2 Insertion: To insert an element, first search for the presence of the element with same key. If such an element is found, insertion fails because duplicates are allowed. When searching is unsuccessful, then insert new element into the last internal node encountered on the search path. To insert an element with key 3 into **Fig 2**, search terminates at second child of left child of the root. It can be inserted into [2,4,6] node results [2,3,4,6] node since this node hold up to 6 elements. The number of disk accesses to do this is 3 in which two accesses for reading root and then its left child and another for writing out modified node after insertion. It can be shown as follows.

To insert an element with key 25 in B-tree of order 7 (Fig 2), the element goes into middle child of the root i.e [20,30,40,50,60,70] but this node is full. When element goes to full node, the overfull node need to be split as follows.

Let the overfull node be $P=[20,25,30,40,50,60,70]$. Let it has m elements and $m+1$ children. It can be denoted as $m, c_0, (e_1, c_1), (e_2, c_2), \dots, (e_m, c_m)$.

where e_i 's indicate elements and c_i 's represent children pointers. The node is split around d where $d=\lfloor m/2 \rfloor$.

The elements to the left remain in P and to the right move into a new node Q but P & Q must contain at least $\lfloor m/2 \rfloor$ children.

The element e_d moved to the parent of P . The format of P and Q are

$P: d-1, c_0, (e_1, c_1), \dots, (e_{d-1}, c_{d-1})$

$Q: m-d, c_d, (e_{d+1}, c_{d+1}), \dots, (e_m, c_m)$.

In this case, the overfull node is $7, 0, (20, 0), (25, 0), (30, 0), (40, 0), (50, 0), (60, 0), (70, 0)$. It can be split around $d=4$ which yields $P=3, 0, (20, 0), (25, 0), (30, 0)$ and $Q=3, 0, (50, 0), (60, 0), (70, 0)$. The $e_4=40$ moved to P 's parent. Here, it is the root. It can be shown as follows. The number of disk accesses required is 5 in which two for searching the proper position in the tree, two for writing out the split nodes and one for writing modified root.

To insert element with key 44 into B – tree of order 3 like Fig 3 (c), the element goes to [35,40] node. Since it is full, the overfull node is [35,40,44] can be represented as $3, 0, (35, 0), (40, 0), (44, 0)$. It can be split around $d=\lfloor 3/2 \rfloor=2$ yields

$P=1, 0, (35, 0)$ and $Q=1, 0, (44, 0)$. The element with key 40 move to P 's parent

$A=[50,60]$. The resulted overfull node be $3, P, (40, Q), (50, C), (60, D)$ where C & D are pointers to the nodes [55] & [60]. The overfull node A be split to create a new node B . The new A & B are $A: 1, P, (40, Q)$ and $B: 1, C, (60, D)$.

Before insertion, root format is R: 2, S, (30,A),(80,T) where S & T are first and third sub trees of the root. After insertion, the overfull node is R: 3, S, (30,A), (50,B),(80,T). This node is split around $d=\lceil 3/2 \rceil=2$ yields R: 1,S,(30,A) and U= 1,B, (80,T). The element 50 moved to R's parent. Since R has no parent, it can be created as new root and that has format 1,R, (50,U). The resulting tree is shown as below.

The total number of disk accesses is 10 in which 3 accesses for reading [30,80],[50,60] and [30,40], six disk accesses for writing out 3 split nodes and one for writing out new root.

When insertion cause s nodes to split, the number of disk accesses is h (reading the nodes on the search path) + 2s (to write out two split parts of each node that is split) + 1 (to write out new root). The total number of disk accesses is $h+2s+1$ which is at most $3h+1$.

3.9.2.3 Deletion: Deletion first divided into 2 cases. (1) The element to be deleted in a node whose children are external nodes. (2) the element to be deleted from a non leaf. case (2) is transformed into case (1) by either largest element in its left neighboring sub tree or smallest element in its right neighboring sub tree.

(i) To delete an element with key 80 in Fig 3 (a), the suitable replacement used is either the largest element in its left sub tree 70 or smallest element 82 in its right sub tree.

(ii) To delete an element with key 80 in Fig 3 (c), the replacing element used is either 70 or 82. If 82 is selected, the problem of deleting 82 from the leaf remains.

The (ii) falls into 2 cases. One is delete an element from a leaf contains more than the minimum number of elements (1 if the leaf is also root and $\lceil m/2 \rceil-1$ if it is not) requires to simply write out modified node.

The deletion from a B-tree of height h is when merging tales at levels h,h-1,...and 3 and getting an element from a nearest sibling at level 2 is $3h$.

Note: when the element size is large relative to the size of a key, the following node structure is used. $s, c_0, (k_1, c_1, p_1), (k_2, c_2, p_2), \dots, (k_s, c_s, p_s)$ where s is the number of elements in the node, k_i 's are element keys, p_i 's are the disk locations of the corresponding elements and c_i 's are children pointers.

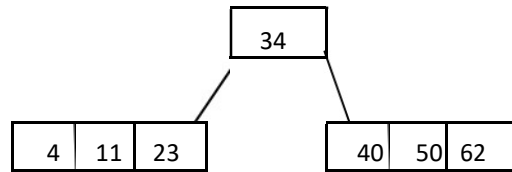
Exercise 1: Draw the B-tree of order 7 resulting from inserting the following keys into an empty tree T: 4,40,23,50,11,34,62,78,66,22,90,59,25,72,64,77,39 & 12.

Step 1: Since it is a B-tree of order 7, the maximum number of elements a node contain is 6.

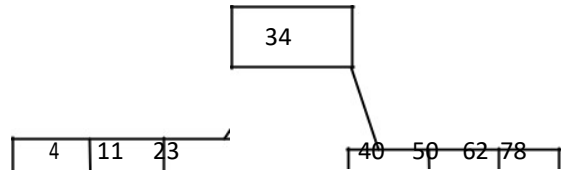
4	11	23	34	40	50
---	----	----	----	----	----

Step 2: Next element to be inserted is 62, but this is full because the maximum number of children that internal node have is 7 and minimum number of children is $4=\lceil 7/2 \rceil$.

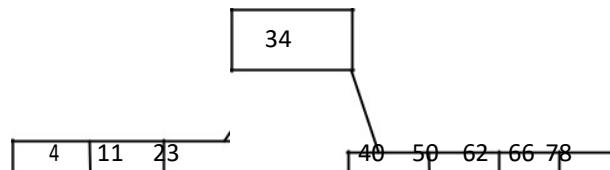
The overfull node is P= [4,11, 23,34, 40,50,62]. It can be split around $e_4=34$. The elements to left are remain in P and to the right in Q. The element e_4 goes to the node parent.



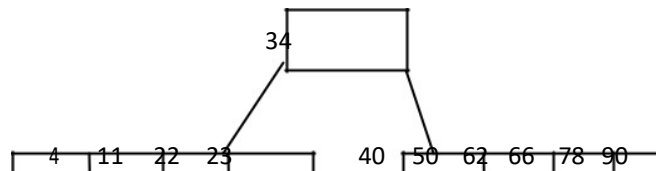
Step 3: The next element 78 goes to root right child.



Step 4: The next element inserted is 66, it goes into root right child.

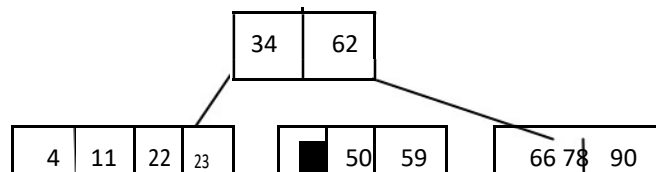


Step 5: The next elements 22 & 90 goes into root left child and root right child respectively. Now, root right child is full. If any element insert into it needs split.

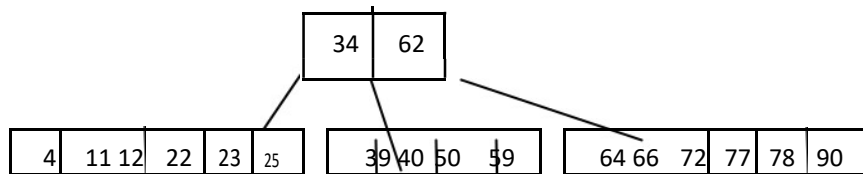


Step 6: The next element 59 goes into root right child and it becomes overfull node. This needs to be split.

Let C= [40,50,59,62,66,78,90]. It can be split around $e_4=62$ leaves C=[40,50,59] and D=[66,78,90]. The element 62 moves to the parent. Now, the root is 34,62]. Its Childs are P,C & D.



Step 7: The elements 25, 72, 64, 77, 39 & 12 are inserted in which 25 & 12 are in root first sub tree, 39 to root middle sub tree and 64, 72 & 77 to root third sub tree.



Huffman coding

is a lossless data compression algorithm. The idea is to assign variable-length codes to input characters, lengths of the assigned codes are based on the frequencies of corresponding characters. The most frequent character gets the smallest code and the least frequent character gets the largest code.

The variable-length codes assigned to input characters are [Prefix Codes](#), means the codes (bit sequences) are assigned in such a way that the code assigned to one character is not the prefix of code assigned to any other character. This is how Huffman Coding makes sure that there is no ambiguity when decoding the generated bitstream.

Let us understand prefix codes with a counter example. Let there be four characters a, b, c and d, and their corresponding variable length codes be 00, 01, 0 and 1. This coding leads to ambiguity because code assigned to c is the prefix of codes assigned to a and b. If the compressed bit stream is 0001, the de-compressed output may be “cccd” or “ccb” or “acd” or “ab”.

See [this](#) for applications of Huffman Coding.

There are mainly two major parts in Huffman Coding

1. Build a Huffman Tree from input characters.
2. Traverse the Huffman Tree and assign codes to characters.

Steps to build Huffman Tree

Input is an array of unique characters along with their frequency of occurrences and output is Huffman Tree.

1. Create a leaf node for each unique character and build a min heap of all leaf nodes (Min Heap is used as a priority queue. The value of frequency field is used to compare two nodes in min heap. Initially, the least frequent character is at root)
2. Extract two nodes with the minimum frequency from the min heap.
3. Create a new internal node with a frequency equal to the sum of the two nodes frequencies. Make the first extracted node as its left child and the other extracted node as its right child. Add this node to the min heap.
4. Repeat steps#2 and #3 until the heap contains only one node. The remaining node is the root node and the tree is complete.

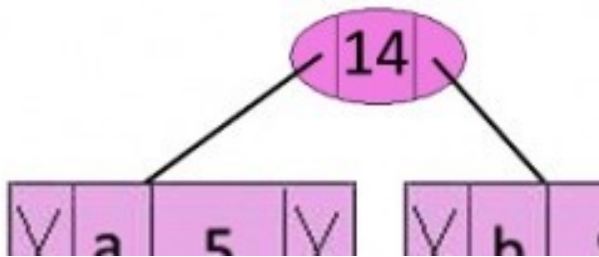
Let us understand the algorithm with an example:

character Frequency

- | | |
|---|----|
| a | 5 |
| b | 9 |
| c | 12 |
| d | 13 |
| e | 16 |

Step 1. Build a min heap that contains 6 nodes where each node represents root of a tree with single node.

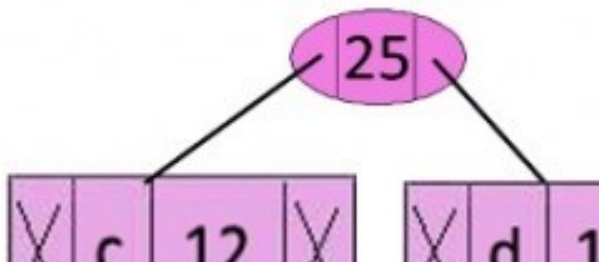
Step 2 Extract two minimum frequency nodes from min heap. Add a new internal node with frequency $5 + 9 = 14$.



Now min heap contains 5 nodes where 4 nodes are roots of trees with single element each, and one heap node is root of tree with 3 elements

character	Frequency
c	12
d	13
Internal Node	14
e	16
f	45

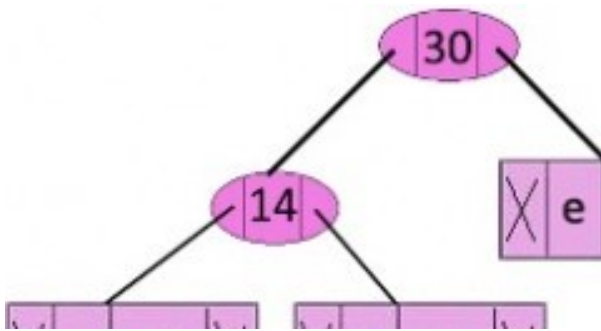
Step 3: Extract two minimum frequency nodes from heap. Add a new internal node with frequency $12 + 13 = 25$



Now min heap contains 4 nodes where 2 nodes are roots of trees with single element each, and two heap nodes are root of tree with more than one nodes

character	Frequency
Internal Node	14
e	16
Internal Node	25
f	45

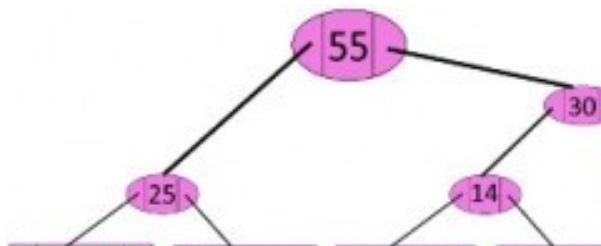
Step 4: Extract two minimum frequency nodes. Add a new internal node with frequency $14 + 16 = 30$



Now min heap contains 3 nodes.

character	Frequency
Internal Node	25
Internal Node	30
f	45

Step 5: Extract two minimum frequency nodes. Add a new internal node with frequency $25 + 30 = 55$

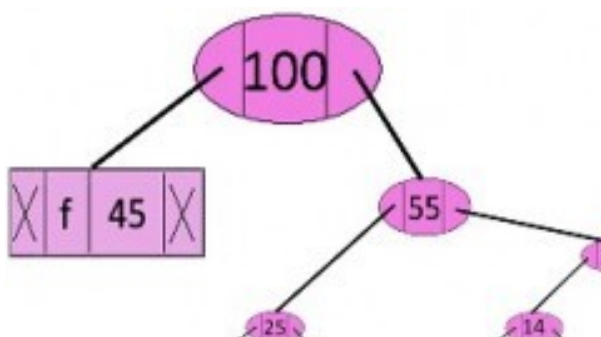


Now min heap contains 2 nodes.

character	Frequency
f	45

Internal Node 55

Step 6: Extract two minimum frequency nodes. Add a new internal node with frequency $45 + 55 = 100$



Now min heap contains only one node.

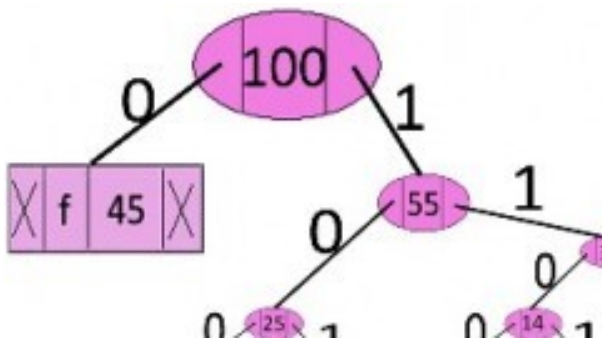
character Frequency

Internal Node 100

Since the heap contains only one node, the algorithm stops here.

Steps to print codes from Huffman Tree:

Traverse the tree formed starting from the root. Maintain an auxiliary array. While moving to the left child, write 0 to the array. While moving to the right child, write 1 to the array. Print the array when a leaf node is encountered.



The codes are as follows:

character code-word

f	0
c	100
d	101
a	1100
b	1101
e	111