

4M26 - Examples Paper 2 - Advanced Data Structures

1.

Briefly answer multiple questions below.

- (i) Explain the general idea of a hash function. If the data you are working with involves string of characters, suggest a scheme that might be used to compute reasonable hash values.

A hash function, $h : U \rightarrow \{0, \dots, m-1\}$, maps every element of a universe of keys $U = \{k_0, k_1, \dots, k_{m-1}\}$ to one element in a set $\{0, \dots, m-1\}$. A desired $h(k)$ is fast to compute and minimises the possibility of collisions between keys. Ideal $h(k)$ could be constructed by rolling a fair m -sided die for each key, k , and using that value to map this key to. The proposed ideal procedure is not feasible if the universe of keys, U , is large as it would not (easily) fit into memory. Hence, other pseudo random functions are often used, e.g. $h(k_i) = k_i \bmod m$ when keys are integer numbers. Note, usually the hash table size, m , is chosen to be a prime number. A good hash function for a string $s = [c_0, c_1, \dots, c_{k-1}]$ where c_i are ASCII characters encoded as integers in range of $\{0, \dots, 255\}$, could be $h(s) = (\sum_{i=0}^{k-1} c_i w_i) \bmod m$. $\{w_i\}$ is a set of randomly selected weights (ideally co-prime). m is also ideally a prime number.

- (ii) Assume your hash table has a space for a maximum of n records in it but only k records are in use at some stage. You can assume that the open addressing scheme with double-hashing is used and the hash function has perfect statistical properties. Estimate the number of table accesses (probes) needed when inserting a new entry into the table.

Let $\alpha = \frac{k}{n}$ be the load factor of the hash table. The probability of a successful insert in one probe is $1 - \alpha$. Probability of a successful insert in exactly t probes is $P(\text{success in } t) = (1 - \alpha)\alpha^{t-1}$. Hence, the expected number of inserts is $E[t] = (1 - \alpha) \sum_{t=0}^{\infty} t\alpha^{t-1} = \frac{1}{1-\alpha}$.

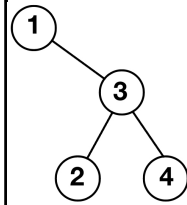
- (iii) Given a binary search tree and a key that is present in it explain an efficient procedure that can be used to update the tree so that the deletion of the given key becomes straight-forward. Indicate expected costs for your procedure. In your answer you can assume that the tree is well balanced. What if the tree is not balanced?

If the node to be deleted has no children or just one child, then deletion is straight-forward and no additional update of the tree is needed. If the node to be deleted has 2 children then swap it with its successor node. Note that the successor could at worst have one child, so after the swap it becomes easy to delete the node of interest.

If the tree is balanced the expected cost is $O(\log n)$. In the worst case, the tree is highly imbalanced (e.g. left side - 1 node, right side $n - 1$ nodes forming a subtree of depth $n - 2$). The node with the key of interest is on the left side of the root node. The successor of the node of interest is at the bottom of the right subtree of the root node. In such case the run time cost of the proposed update procedure is $O(n)$ as it takes $O(n)$ steps to find the successor.

- (iv) Suppose that the search for key, k , in a binary search tree ends up in a leaf. Consider three sets: A - the keys to the left of the search path; B - the keys on the search path; and C - the keys to the right of the search path. Is the following claim correct: *any three keys $a \in A$, $b \in B$, and $c \in C$ must satisfy $a \leq b \leq c$* ? Justify your answer.

Contradiction by example.



Lets assume that we search for 4. Then $A = \{2\}$, $B = \{1, 3, 4\}$ and $C = \emptyset$. Since $2 > 1$ the claim is not true.

(v) Given a binary search tree and a key that is present in it, describe a procedure that rearranges the tree so that the specified key is in the root of the updated tree. Is this always possible?

Consider a *left rotation* operation, described in the lecture slides for red black trees. It takes a tree with nodes M , N such that M is the left child of N and tree T_3 is the right child of N . Also trees T_1 and T_2 are correspondingly left and right children of M . Such a tree is transformed into a new tree with the *left rotation*: M has T_1 and N as its left and right children correspondingly. Where N has T_2 and T_3 as its left and right children. The *right rotation* performs a similar procedure to bring the right child of a particular node up. The ordering property of the binary search tree is preserved by both operations, but moves the relevant node (left child for the left rotation and right child for the right rotation) up one level. By applying it repeatedly, the problem can be solved with a cost proportional to the distance upwards that the node of interest must move.

(vi) What is a *red-black tree*? What are the key properties of red-black trees and what problem do they solve?

Red-black trees are binary search trees that are managed in such a way that they remain tolerably balanced, thereby avoiding the key problem that simple binary search trees suffer from. In red-black trees each node can be tagged "red" or "black". The red node can only appear as a direct descendent o a black node.

To be precise, red-black trees are binary search trees that satisfy 5 properties:

1. Every node is *red* or *black*.
2. The root node is *black*.
3. Every leaf node (nil node) is *black*.
4. If a node is *red*, both of its children are *black*.
5. Starting from any node, all simple paths down to leaf nodes hold the same number of *black* nodes.

All the aforementioned properties are maintained when performing same operations as performed on binary search trees: *search*, *insert*, *delete*, etc.

(vii) Suppose that a priority queue is to be implemented using red-black trees. How can the smallest item in the queue be found? Starting from an empty queue how much work is required to build a queue with n items in it?

The smallest item is the leftmost leaf of the red-black tree and can be reached in $O(\log n)$ steps. Building a tree with n nodes takes $\sum_{k=1}^n O(\log k)$, steps where k is the size of the queue prior each insertion. Hence a total of $O(\log n!)$ steps are required, which is $O(n \log n)$.

(viii) Give an upper and a lower bound for the number of nodes in B-tree of height, h , and degree, t . Hence if a B-tree holds n values find an upper bound on its height.

A B-Tree with degree, t , has nodes with number of children ranging from t to $2t$. Hence the lower bound in terms of number of nodes is t^h and upper bound is $(2t)^h$. Therefore the height of a B-tree of n nodes is at most $\log_t n$.

(x) Explain how search works in a B-tree. In your answer you can ignore disk read and write operations.

Search starts at the root node. For each node, the key which we are searching for is compared with the keys stored in the node (in an increasing order) until the key which is equal or larger is found. In the first case the current node and the index of the relevant key is returned. In the second case the search is continued in the subtree of the child which stores values larger than a key smaller than the key that we are searching for and smaller than the key that was just found. An example python code is provided below. The runtime of the search algorithm described above is $O(t \log_t n)$ since in the worst case one may need to scan every key in each node on a path in a B-tree the maximum length of which is $\log_t n$.

```
In [1]: def search(self, u, key): # u is a node
        # linear scan to find index of key
        i = 0
        while i < len(u.keys) and key > u.keys[i]:
            i += 1
        if i < len(u.keys) and key == u.keys[i]:
            return (u, i)
        if u.is_leaf:
            return None

        return self.search(u.children[i], key)
```

2.

Write the function, **rolling_hash(S)**, which:

- (i) takes a two element list, S , which contains two strings of length n and m correspondingly and;
- (ii) finds the index of the first location in the string, $S[1]$, starting from which the string, $S[0]$, appears in string, $S[1]$, in full. If the first string is not a substring of the second string, your function should output, -1 .

The expected runtime complexity of your method should be $O(m + n)$ when a substring can not be found and the expected number of collisions in the hash table is constant for a typical input.

Examples:

Input: ["abc","abcdefgh"]

Output: 0

Input: ["fgh","abcdefgh"]

Output: 5

Constraints:

- $1 \leq n \leq 1000$.
- $1 \leq m \leq 1000$.

Code:

If one hashes a string, s , of length, n , as $h(s) = (\sum_i s_{n-i-1} \cdot p^i) \bmod m$, then recomputing a hash for a string $s' = s[1:] + c$, where $+$ is a concatenation operation of a string, $s[1:]$, and a character, c , can be performed in constant time.

```
In [2]: def geometric_hash(s,m,p):
        val=0
```

```

p_loop = 1
n = len(s)
for i in range(n):
    val= (val+ord(s[n-i-1])* p_loop) % m
    p_loop*=p_loop
return val

def rolling_hash(S):

    #Write your code here

    pattern = S[0]
    data = S[1]
    m = len(pattern)
    n = len(data)

    if m>n:
        return -1

    p=7699
    p_m = pow(p,m) % m

    hpattern = geometric_hash(pattern,m,p)

    hs = geometric_hash(data[:m],m,p)

    for i in range(m,n):
        if hs == hpattern and data[(i-m):i] == pattern:
            return i-m
        hs = (hs * p - ord(data[i-m]) * p_m + ord(data[i]))%m

    if hs == hpattern and data[(n-m):n] == pattern:
        return n-m
    return -1

```

Tests:

Run example test case 1:

```

In [3]: input_value = ["abc", "abcdefgh"]
        print (rolling_hash(input_value))

```

0

Run example test case 2:

```

In [4]: input_value = ["fgh", "abcdefgh"]
        print (rolling_hash(input_value))

```

5

3.

Write functions, **inorder**(node), **preorder**(node), **postorder**(node), each of which traverses a binary search tree in the corresponding order.

You are provided with the function, **bst_manipulate**(insert_keys), which takes a list of keys as an input, creates a binary search tree (BST), using *BinarySearchTree* class and calls the aforementioned BST traversal functions to be implemented. The list of three lists containing keys in the order of inorder, preorder and postorder traversal of the BST is returned as an output.

Examples:

Input: [2, 1, 3]

Output: [[1, 2, 3],[2, 1, 3],[1, 3, 2]]

Input: [5, 3, 2, 7, 1, 8, 9, 12]

Output: [[1, 2, 3, 5, 7, 8, 9, 12],[5, 3, 2, 1, 7, 8, 9, 12],[1, 2, 3, 12, 9, 8, 7, 5]]

Constraints:

- Keys are assumed to be unique, integer values.
- $1 \leq \text{len}(\text{insert_keys}) \leq 100$.

Code:

In [5]:

```
from utils.binary_search_tree import BinarySearchTree, Node, insert
from utils.binary_search_tree import inorder as inorder_orig
from utils.binary_search_tree import preorder as preorder_orig
from utils.binary_search_tree import postorder as postorder_orig

def inorder(node):
    #Write your code here

    return inorder_orig(node)

def preorder(node):
    #Write your code here

    return preorder_orig(node)

def postorder(node):
    #Write your code here

    return postorder_orig(node)

def bst_manipulate(insert_keys):

    bst = BinarySearchTree()

    node_list = [Node(key) for key in insert_keys]

    for node in node_list:
        insert(bst, node)

    return [inorder(bst.root),preorder(bst.root),postorder(bst.root)]
```

Tests:

Run example test case 1:

In [6]:

```
input_value = [5, 3, 2, 7, 1, 8, 9, 12]
print (bst_manipulate(input_value))
```

[[1, 2, 3, 5, 7, 8, 9, 12], [5, 3, 2, 1, 7, 8, 9, 12], [1, 2, 3, 12, 9, 8, 7, 5]]

Run example test case 2:

In [7]:

```
input_value = [2, 1, 3]
print (bst_manipulate(input_value))
```

[[1, 2, 3], [2, 1, 3], [1, 3, 2]]

4.

Write the function, `delete(bst, node)`, which deletes a node, `node`, from a binary search tree, `bst`.

You are provided with the function, `bst_delete(delete_keys)`, which takes a list of keys to be deleted as an input, creates a binary search tree (BST), using a provided `BinarySearchTree` class, and then calls the `delete` function for nodes with keys in `delete_keys` list.

Note, this question can be tackled at varying difficulty. In the easiest case, access to functions `minimum(node)`, `shift_nodes(bst, old_node, new_node)` is allowed. The former function finds the node with the smallest key in the subtree of node, `node`. The latter function, shifts the nodes from the subtree at `new_node` to the position of `old_node`.

In the hardest case, one should attempt to implement the whole delete procedure, including the aforementioned functions.

Examples:

In all test examples it is assumed that the tree is first built by inserting nodes with the following keys: `[5, 3, 2, 7, 1, 8, 9, 12]`.

Input: `[5]`

Output: `[7, 3, 2, 1, 8, 9, 12]`

Input: `[7, 1, 8]`

Output: `[5, 3, 2, 9, 12]`

Constraints:

- Keys are assumed to be unique, integer values.
- $1 \leq \text{len}(\text{insert_keys}) \leq 100$.
- Each key that is about to be deleted is guaranteed to exist in the tree.

Code:

```
In [8]: from utils.binary_search_tree import BinarySearchTree, Node, insert, preorder, search
        from utils.binary_search_tree import minimum, shift_nodes

        def delete(bst, node):

            #Write your code here

            if not node.right: # node to be deleted has no right child
                shift_nodes(bst, node, node.right) # node to be deleted has no left child
            elif not node.left:
                shift_nodes(bst, node, node.right)
            else: # node to be deleted has both a left and right child
                successor = minimum(node.right)
                if successor != node.right:
                    shift_nodes(bst, successor, successor.right)
                    successor.right = node.right
                    successor.right.parent = successor
                shift_nodes(bst, node, successor)
                successor.left = node.left
                successor.left.parent = successor

        def bst_delete(delete_keys):

            bst = BinarySearchTree()
```

```

insert_keys = [5, 3, 2, 7, 1, 8, 9, 12]
node_list = [Node(key) for key in insert_keys]

for node in node_list:
    insert(bst, node)

for key in delete_keys:
    node = search(key, bst.root)
    delete(bst, node)

return preorder(bst.root)

```

Tests:

Run example test case 1:

```

In [9]: input_value = [5]
        print (bst_delete(input_value))

```

[7, 3, 2, 1, 8, 9, 12]

Run example test case 2:

```

In [10]: input_value = [7,1,8]
         print (bst_delete(input_value))

```

[5, 3, 2, 9, 12]

5.

Below you can find a partial code (*IntervalRedBlackTree* class) to implement an interval tree by building upon a red-black tree (*RedBlackTree* class).

An interval is represented by a list of two numbers (e.g. [*low*,*high*]) indicating the lower and higher end of the interval correspondingly. An interval tree is a binary search tree in which all intervals of all the nodes on the left child subtree of a particular node have their lower value lower than the *low* value of this node and all intervals of the nodes on the subtree of the right child have their lower value higher than the *low* value of the interval of this particular node. I.e. it is similar to a binary search tree where keys are represented by lower ends of intervals. However, additional information is stored in the node attributes (see *IntervalNode* class).

An important operation on interval trees is finding a node which overlaps with an interval of interest. Write the function, **search_interval**(*self*,*interval*), which implements the aforementioned functionality for a *IntervalRedBlackTree* class. Note, your implementation may need a function, **update_max_value**(*self*,*node*), which for a particular node and its children updates the attribute, *max_val*, which stores to the maximum *high* value of all intervals in the subtree of that particular node.

Note, you are provided with the function, **manipulate_interval_tree**(*interval*), which creates an interval tree, calls **update_max_value** function for the root node (runtime $O(n)$, n - number of intervals in a tree) and returns the output of the **search_interval** function.

Examples:

Input: [11,13]

Output: Node(key/low = 13, high=14 color=red, max_val=14), parent=9, left=Nil, right=Nil

Input: [-7,-9]

Output: Node(Nil, color=black)

Constraints:

- $1 \leq n \leq 100$.

Code:

In [11]:

```
#A straight-forward extension of the Node of a RedBlackTree to store the information
#the interval.

from utils.red_black_tree import RedBlackTree, Node, Nil

class IntervalNode(Node):

    def __init__(self, low, high, parent=None, left=None, right=None, color=None, max_val=None):
        self.high = high #High end of an interval
        self.max_val = max_val, #The largest high end value of all intervals of the subtree
        super().__init__(low, parent = parent, left = left, right = right, color = color)

    def __repr__(self):
        summary = f"Node(key/low = {self.key}, high={self.high} color={self.color}, max_val={self.max_val}"

        if self.parent:
            summary += f", parent={self.parent.key}"
        if self.left:
            summary += f", left={self.left.key}"
        if self.right:
            summary += f", right={self.right.key}"
        return summary

class IntervalRedBlackTree(RedBlackTree):

    def __init__(self):
        super().__init__()

    def update_max_value(self, node):
        #Write your code here

        if not (node is self.nil):
            self.update_max_value(node.left)
            self.update_max_value(node.right)
            node.max_val = node.high
            if not (node.left is self.nil) and node.left.max_val > node.max_val:
                node.max_val = node.left.max_val
            if not (node.right is self.nil) and node.right.max_val > node.max_val:
                node.max_val = node.right.max_val

    def search_interval(self, interval):
        #Write your code here

        x = self.root
        while not (x is self.nil) and (interval[1] < x.key or interval[0] > x.high):
            if not (x.left is self.nil) and x.left.max_val >= interval[0]:
                x = x.left
            else:
                x = x.right
        return x

    def manipulate_interval_tree(interval):

        rbt = IntervalRedBlackTree()

        interval_list = [[5,6], [3,4], [2,3], [7,8], [1,2], [8,9], [9,10], [13,14]]
```



```

node_list = [IntervalNode(interval[0],interval[1]) for interval in interval_list]

for node in node_list:
    rbt.insert(node)

rbt.update_max_value(rbt.root)

return rbt.search_interval(interval)

```

Tests:

Run example test case 1:

```

In [12]: input_value = [11,13]
         print (manipulate_interval_tree(input_value))

```

```

Node(key/low = 13, high=14 color=red, max_val=14), parent=9, left=Nil, right=Nil

```

Run example test case 2:

```

In [13]: input_value = [-7,-9]
         print (manipulate_interval_tree(input_value))

```

```

Node(Nil, color=black)

```

6.

A B-tree of minimum degree, t , is a data structure with the property that every node may contain at most $2t - 1$ keys. For this question, we will assume that all keys are distinct, that all nodes contain keys but not values, and that all nodes of the B-tree fit into memory (no disk read or writes are required).

Implement a function to insert a given key, key , in a B-tree by completing the **insert_key**(key) function below. For a B-tree storing n keys, your solution should run in $O(t \log_t n)$ time.

Note that the function, **insert_key**, is implemented for a class, *Btree_new*, which is derived from an original *Btree* implementation provided with the lectures. Additional code is also provided in function, **btree_manipulate**(key_list), which helps to initialise a B-tree, calls your function to insert a list of keys, visualises it, and then deletes a multiple keys in a sequence from a Btree. A list key_list containing two lists of keys to insert and keys to delete is provided as an input to **btree_manipulate**. The minimum and maximum keys left after insertion and deletion steps are performed should be returned by **btree_manipulate**.

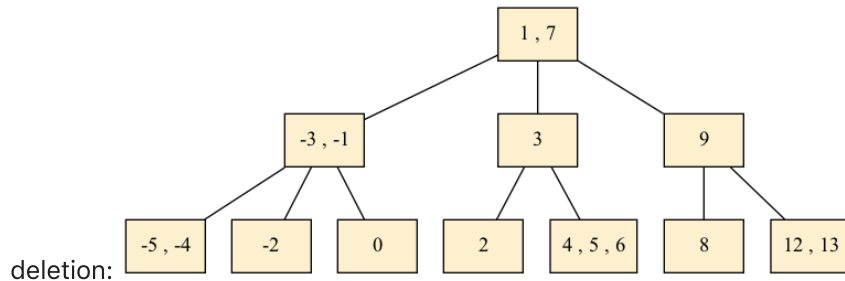
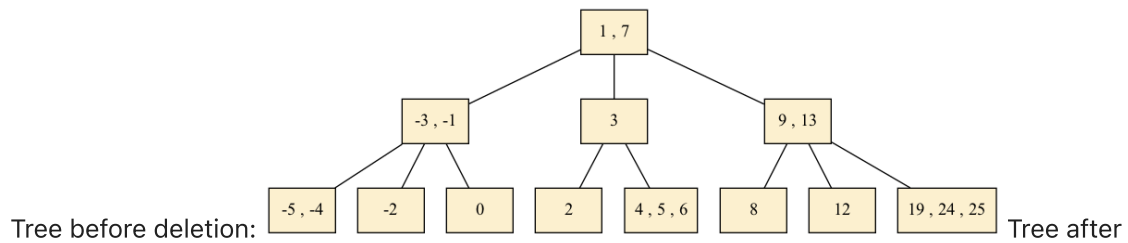
Remark. Insertion operation in a Btree is not examinable, however, it is a very useful exercise to implement this operation as it helps to understand the inner workings of the Btree. If you find that the implementation is taking too long, you can attempt implementing a search function instead.

Examples:

Input: [[5, 3, 2, 7, 1, 8, 9, 12, 13, 4, 0, 6, -1, 19, 24, 25, -2, -3, -4, -5], [25,24,19]]

Output: Minimum key: -5, maximum key: 24

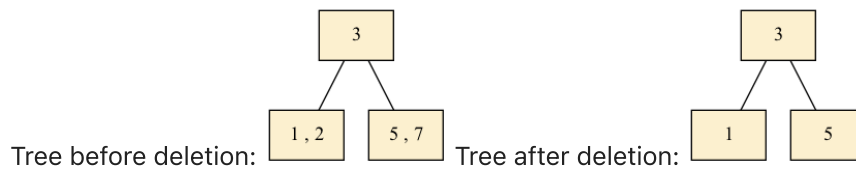
Explanation:



Input: `[[5, 3, 2, 7, 1], [2, 7]]`

Output: Minimum key: 1, maximum key: 5

Explanation:



Constraints:

- $1 \leq n \leq 100$.
- Each key that is about to be deleted is guaranteed to exist in the B-tree.

Code:

In [14]:

```
from utils.btree import Btree, Node
from pathlib import Path
from matplotlib import pyplot as plt
import cv2

class Btree_new(Btree):
    def __init__(self, t):
        super().__init__(t, root=Node(is_leaf=True), verbose=False)

    def insert_key(self, key):
        # Write your code here
        self.insert(key)

def btree_manipulate(key_list):
    insert_keys, keys_to_delete = key_list[0], key_list[1]

    btree = Btree_new(t=2)

    for key in insert_keys:
        btree.insert_key(key)

    btree.viz_btree(dest_path=Path("btree.png"), refresh=True)

    print('Before deletion:')
    plt.imshow(cv2.imread("btree.png")[:, :, ::-1])
    plt.show()

    print("Keys to be deleted:")
```

```

print(keys_to_delete)

for key in keys_to_delete:
    btree.delete(btree.root, key)

btree.viz_btree(dest_path=Path("btree-after-deletions.png"), refresh=True)

print('After deletion:')
plt.imshow(cv2.imread("btree-after-deletions.png")[:, :, ::-1])
plt.show()

return "Minimum key: "+str(btree.minimum(btree.root))+', maximum key: '+str(btree

```

Tests:

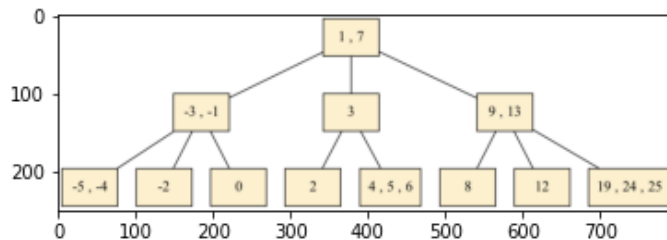
Run example test case 1:

```

In [15]: input_value = [[5, 3, 2, 7, 1, 8, 9, 12, 13, 4, 0, 6, -1, 19, 24, 25, -2, -3, -4, -5]
print (btree_manipulate(input_value))

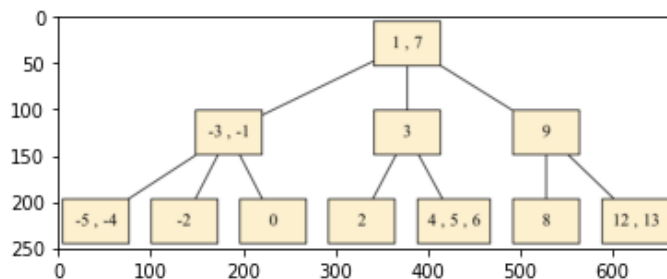
```

Saving visualisation to btree.png
Before deletion:



Keys to be deleted:
[25, 24, 19]

Saving visualisation to btree-after-deletions.png
After deletion:



Minimum key: -5, maximum key: 13

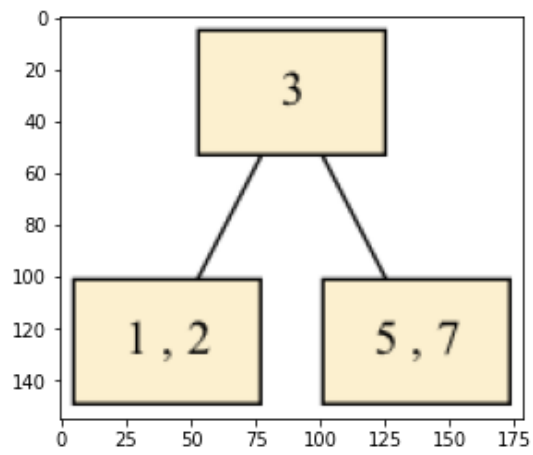
Run example test case 2:

```

In [16]: input_value = [[5, 3, 2, 7, 1],[2,7]]
print (btree_manipulate(input_value))

```

Saving visualisation to btree.png
Before deletion:

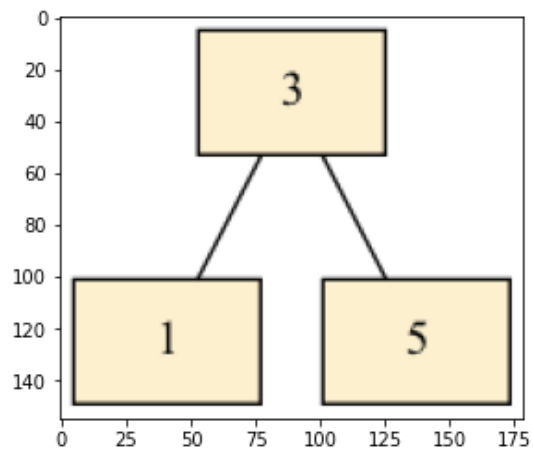


Keys to be deleted:

[2, 7]

Saving visualisation to btree-after-deletions.png

After deletion:



Minimum key: 1, maximum key: 5