

Algorithms: Order Statistics Red Black Tree

2017-18570 Sungchan Yi

Dept. of Computer Science and Engineering

May 15th, 2019

Contents

| | | |
|----------|------------------------|----------|
| 1 | Introduction | 1 |
| 2 | Implementation | 2 |
| 2.1 | Checker | 2 |
| 3 | Experiments | 4 |
| 3.1 | Test Case #1 | 6 |
| 3.2 | Test Case #2 | 7 |
| 4 | Conclusion | 9 |

1 Introduction

In this homework, we implement an **order statistic red black tree** with C++ language using data structure augmentation. Compared to normal red black trees, order statistics tree will contain an additional field *size* in each node v , which keeps track of the number of nodes in the subtree with root v . Additionally, we implement a checker program for the correctness for the implemented data structure and operations. The implemented tree should support the following operations:

- OS-INSERT(T, x): Insert x into T . If x already exists in T , return 0. Otherwise, return x .
- OS-DELETE(T, x): Delete x from T . If x is not in T , return 0. Otherwise, return x .
- OS-SELECT(T, i): Select the i -th smallest key from T and return it. If i is greater than the number of keys currently in T , return 0.

- OS-RANK(T, x): Return the rank of x if x is in T . Otherwise, return 0.

The following are the restrictions on given inputs.

- Each operation will be denoted as I, D, S, R respectively.
- $1 \leq x \leq 9999$.
- Input will be given in the form “[Character representing the operation] [x]”.

2 Implementation

The basic operations required for the tree was implemented as the same way in the textbook, so we omit it here.

2.1 Checker

To check the correctness, we implement a checker program. The basic idea is to use a `vector<int> vec` to keep track of the set of integers given as input.

- `check_insert(int x)` will traverse `vec` and check if `x` is already in the vector. If it is, immediately return 0. If it is not found, `push_back x` and return `x`.

```

1  int check_insert(int x) {
2      for(auto v : vec) {
3          if(v == x) return 0;
4      }
5      vec.push_back(x);
6      return x;
7  }
```

- `check_delete(int x)` will traverse `vec` and check if `x` is already in the vector. If it is, swap it with the last element of `vec`, `pop_back`, and return `x`. If it is not found, return 0.

```

1  int check_delete(int x) {
2      for(int i = 0; i < vec.size(); ++i) {
3          if(x == vec[i]) {
4              int tmp = vec[vec.size() - 1];
5              vec[vec.size() - 1] = x;
6              vec[i] = tmp;
7              vec.pop_back();
8              return x;
9          }
10     }
11     return 0;
12 }
```

- `check_select(int x)` will first check if `x` is less than the size of `vec`. If it is, sort the vector and return `vec[x - 1]`. Else, return 0.

```

1 int check_select(int x) {
2     if(x <= vec.size()) {
3         sort(vec.begin(), vec.end());
4         return vec[x - 1];
5     } else return 0;
6 }

```

- `check_rank(int x)` will sort the vector and then traverse the list. If `x` is found, return the (index of `x`) + 1. If it is not found, return 0.

```

1 int check_rank(int x) {
2     sort(vec.begin(), vec.end());
3     for(int i = 0; i < vec.size(); ++i) {
4         if(vec[i] == x) return i + 1;
5     }
6     return 0;
7 }

```

And here is the checker code. Interpret `opt_seq` and perform `check_[op]`, then compare it with the given `out_seq`. If the result doesn't match in any point in time, return `false`.

```

1 bool check(int opt_seq[], int in_seq[], int out_seq[], int n) {
2     string str = "";
3     for(int i = 0; i < n; ++i) {
4         if(opt_seq[i] == 0) {
5             if(check_insert(in_seq[i]) != out_seq[i]) {
6                 printf("%d: WA on I %d\n", i, in_seq[i]);
7                 return false;
8             }
9         } else if(opt_seq[i] == 1) {
10            if(check_delete(in_seq[i]) != out_seq[i]) {
11                printf("%d: WA on D %d\n", i, in_seq[i]);
12                return false;
13            }
14        } else if(opt_seq[i] == 2) {
15            if(check_select(in_seq[i]) != out_seq[i]) {
16                printf("%d: WA on S %d\n", i, in_seq[i]);
17                return false;
18            }
19        } else if(opt_seq[i] == 3) {
20            if(check_rank(in_seq[i]) != out_seq[i]) {
21                printf("%d: WA on R %d\n", i, in_seq[i]);
22                return false;
23            }
24        }
25    }
26 }

```

```

25     }
26     return true;
27 }

```

3 Experiments

Here is how the test was done.

- Test environment

- OS: Ubuntu 18.04.2 LTS
- CPU: Intel Core i5-6200U CPU @ 2.30GHz × 4
- RAM: 8GB

- Test Case Generator (Python)

We keep the numbers small so that the numbers generated in the test will overlap easily, which will cause more operations to return some value other than 0.

```

1  from random import randint
2
3  # Select random element from list
4  def rndselect(list):
5      idx = randint(0, len(list) - 1)
6      return list[idx]
7
8  # Number of Test Cases
9  N = 100000
10
11 # Operation set
12 op_set = ['I', 'D', 'R', 'S']
13
14 # Numbers
15 num = [str(i) for i in range(1, 100)]
16
17 # Create input text file
18 f = open('tests/input', "w")
19
20 # Write to test file
21 f.write(str(N) + "\n")
22 for i in range(N):
23     f.write(rndselect(op_set))
24     f.write(' ')
25     f.write(rndselect(num))
26     f.write('\n')
27 f.close()
28

```

- Checker Code

This is the checker code I implemented, to check the correctness of my implementation. An helper function `printSideways` was used to examine the tree structure easily.

```

1      #include <cstdio>
2      #include "hw2_cpp.h"
3      #define N 101010
4
5      extern Tree tree; // Tree defined in hw2_cpp.h
6
7      int opt_seq[N], in_seq[N], out_seq[N];
8
9      int main() {
10         string str = "";
11
12         freopen("./tests/input", "r", stdin);
13         string op;
14         int x, t;
15         cin >> t;
16         for(int i = 0; i < t; ++i) {
17             cin >> op >> x;
18             in_seq[i] = x;
19             printf("%d. %c %d : ", i + 1, op[0], x);
20             if(op[0] == 'I') {
21                 opt_seq[i] = 0;
22                 out_seq[i] = os_insert(x);
23             } else if(op[0] == 'D') {
24                 opt_seq[i] = 1;
25                 out_seq[i] = os_delete(x);
26             } else if(op[0] == 'S') {
27                 opt_seq[i] = 2;
28                 out_seq[i] = os_select(x);
29             } else if(op[0] == 'R') {
30                 opt_seq[i] = 3;
31                 out_seq[i] = os_rank(x);
32             }
33             printf("%d\n", out_seq[i]);
34             tree.printSideways(tree.getRoot(), str);
35         }
36
37         puts("Checking...");
38         if(check(opt_seq, in_seq, out_seq, t)) puts("correct");
39         else puts("incorrect");
40     }

```

- `printSideways` function

Prints the *data* and *size* of each node, with colors embedded.

```
1 void Tree::printSideways(Node* node, string& indent) {
2     string str = indent + "          ";
3     if(node != NIL) {
4         printSideways(node -> right, str);
5         if(node -> color == RED) {
6             cout << indent;
7             printf("\033[1;31m(%d / %d)\033[0m\n", node -> data,
8 node -> size);
9         } else {
10            cout << indent;
11            printf("\033[1;30m(%d / %d)\033[0m\n", node -> data,
12 node -> size);
13        }
14        printSideways(node -> left, str);
15    }
16 }
```

- Test case generator code was saved to `tester.py`, the checker code was saved to `os_test.cpp`.

- Generate Test: `python3 tester.py`
- Compile: `g++ os_test.cpp -o test`
- Run: `./test`

3.1 Test Case #1

I tried a small test case first, with small numbers.

10 I 3 I 13 I 2 S 1 R 6 R 5 D 10 I 17 I 7 I 8

Here is the structure of the tree after each operation. Each node is represented by (*key/size*). The leftmost node is the root, and by each column, the nodes in the next level are shown.

| | |
|--------------|--------------|
| 1. I 3 : 3 | 6. R 5 : 0 |
| (3 / 1) | (13 / 1) |
| | (3 / 3) |
| | (2 / 1) |
| 2. I 13 : 13 | 7. D 10 : 0 |
| (13 / 1) | (13 / 1) |
| (3 / 2) | (3 / 3) |
| | (2 / 1) |
| 3. I 2 : 2 | 8. I 17 : 17 |
| (13 / 1) | (17 / 1) |
| (3 / 3) | (13 / 2) |
| (2 / 1) | (3 / 4) |
| | (2 / 1) |
| 4. S 1 : 2 | 9. I 7 : 7 |
| (13 / 1) | (17 / 1) |
| (3 / 3) | (13 / 3) |
| (2 / 1) | (7 / 1) |
| | (3 / 5) |
| | (2 / 1) |
| 5. R 6 : 0 | 10. I 8 : 8 |
| (13 / 1) | (17 / 1) |
| (3 / 3) | (13 / 4) |
| (2 / 1) | (8 / 1) |
| | (7 / 2) |
| | (3 / 6) |
| | (2 / 1) |
| | Checking... |
| | correct |

Figure 1: Result of Test 1

3.2 Test Case #2

For the second case, I tried a case with size 30.

30 I 9 I 1 S 12 R 19 R 14 S 6 S 11 R 9 I 3 D 10 S 5 S 17 D 13 I 19
 I 19 S 6 R 3 S 7 I 1 S 10 R 8 I 3 I 10 I 2 I 11 I 15 I 13 I 16 R 15 D
 6

| | | |
|---|---|--|
| 1. I 9 : 9 (9 / 1) | 7. S 11 : 0 (9 / 2) (1 / 1) | 13. D 13 : 0 (9 / 1) (3 / 3) (1 / 1) |
| 2. I 1 : 1 (9 / 2) (1 / 1) | 8. R 9 : 2 (9 / 2) (1 / 1) | 14. I 19 : 19 (19 / 1) (9 / 2) (3 / 4) (1 / 1) |
| 3. S 12 : 0 (9 / 2) (1 / 1) | 9. I 3 : 3 (9 / 1) (3 / 3) (1 / 1) | 15. I 19 : 0 (19 / 1) (9 / 2) (3 / 4) (1 / 1) |
| 4. R 19 : 0 (9 / 2) (1 / 1) | 10. D 10 : 0 (9 / 1) (3 / 3) (1 / 1) | 16. S 6 : 0 (19 / 1) (9 / 2) (3 / 4) (1 / 1) |
| 5. R 14 : 0 (9 / 2) (1 / 1) | 11. S 5 : 0 (9 / 1) (3 / 3) (1 / 1) | 17. R 3 : 2 (19 / 1) (9 / 2) (3 / 4) (1 / 1) |
| 6. S 6 : 0 (9 / 2) (1 / 1) | 12. S 17 : 0 (9 / 1) (3 / 3) (1 / 1) | 18. S 7 : 0 (19 / 1) (9 / 2) (3 / 4) (1 / 1) |
| 19. I 1 : 0 (19 / 1) (9 / 2) (3 / 4) (1 / 1) | 25. I 11 : 11 (19 / 2) (11 / 1) (10 / 4) (9 / 1) (3 / 7) (2 / 1) (1 / 2) | 28. I 16 : 16 (19 / 2) (16 / 1) (15 / 5) (13 / 1) (11 / 2) (10 / 10) (9 / 1) (3 / 4) (2 / 1) (1 / 2) |
| 20. S 10 : 0 (19 / 1) (9 / 2) (3 / 4) (1 / 1) | 26. I 15 : 15 (19 / 1) (15 / 3) (11 / 1) (10 / 5) (9 / 1) (3 / 8) (2 / 1) (1 / 2) | 29. R 15 : 8 (19 / 2) (16 / 1) (15 / 5) (13 / 1) (11 / 2) (10 / 10) (9 / 1) (3 / 4) (2 / 1) (1 / 2) |
| 21. R 8 : 0 (19 / 1) (9 / 2) (3 / 4) (1 / 1) | 27. I 13 : 13 (19 / 1) (15 / 4) (13 / 1) (11 / 2) (10 / 9) (9 / 1) (3 / 4) (2 / 1) (1 / 2) | 30. D 6 : 0 (19 / 2) (16 / 1) (15 / 5) (13 / 1) (11 / 2) (10 / 10) (9 / 1) (3 / 4) (2 / 1) (1 / 2) |
| 22. I 3 : 0 (19 / 1) (9 / 2) (3 / 4) (1 / 1) | | Checking... correct |
| 23. I 10 : 10 (19 / 1) (10 / 3) (9 / 1) (3 / 5) (1 / 1) | | |
| 24. I 2 : 2 (19 / 1) (10 / 3) (9 / 1) (3 / 6) (2 / 1) (1 / 2) | | |

Figure 2: Result of Test 2

I ran many other cases, up to 100000 operations. They were all correct.

4 Conclusion

In this homework, I learned how to augment data structures to contain additional information, that will help me retrieve particular information. With the augmented data structure, the additional field could be calculated in $O(1)$ time, since it only needs the information from left/right child. In total, **rank**, **select** takes $O(\log n)$ time with the order statistics tree, since the height of a red black tree is at most $O(\log n)$.