

Section 04: Dictionaries

Review Problems

1. AVL Trees

Draw an AVL Tree as each of the following keys are added in the order given. Show intermediate steps.

(a)

{“penguin”, “stork”, “cat”, “fowl”, “moth”, “badger”, “otter”, “shrew”, “lion”, “raven”, “bat”}

(b)

{6, 43, 7, 42, 59, 63, 11, 21, 56, 54, 27, 20, 36}

(c)

{“indigo”, “fuchsia”, “pink”, “goldenrod”, “violet”, “khaki”, “red”, “orange”, “maroon”, “crimson”, “green”, “mauve”}

2. True or false?

(a) An insertion in an AVL tree with n nodes requires $\Theta(\log(n))$ rotations.

(b) A set of numbers are inserted into an empty BST in sorted order and inserted into an empty AVL tree in random order. Listing all elements in sorted order from the BST is $\mathcal{O}(n)$, while listing them in sorted order from the AVL tree is $\mathcal{O}(\log(n))$.

(c) If items are inserted into an empty BST in sorted order, then the BST's `get()` is just as asymptotically efficient as an AVL tree whose elements were inserted in unsorted order.

(d) An AVL tree will always do a maximum of two rotations in an insert.

3. Big- \mathcal{O}

Write down the big- \mathcal{O} for each of the following:

(a) Insert and find in a BST.

(b) Insert and find in an AVL tree.

(c) Finding the minimum value in an AVL tree containing n elements.

(d) Finding the k -th largest item in an AVL tree containing n elements.

(e) Listing elements of an AVL tree in sorted order

Section Problems

4. Hash table insertion

For each problem, insert the given elements into the described hash table. Do not worry about resizing the internal array.

- (a) Suppose we have a hash table that uses separate chaining and has an internal capacity of 12. Assume that each bucket is a linked list where new elements are added to the front of the list.

Insert the following elements in the EXACT order given using the hash function $h(x) = 4x$:

0, 4, 7, 1, 2, 3, 6, 11, 16

- (b) Suppose we have a hash table that uses linear probing and has an internal capacity of 13.

Insert the following elements in the EXACT order given using the hash function $h(x) = 3x$:

2, 4, 6, 7, 15, 13, 19

- (c) Suppose we have a hash table that uses quadratic probing and has an internal capacity of 10.

Insert the following elements in the EXACT order given using the hash function $h(x) = x$:

0, 1, 2, 5, 15, 25, 35

5. Analyzing dictionaries

- (a) What are the constraints on the data types you can store in an AVL tree? When is an AVL tree preferred over another dictionary implementation, such as a HashMap?
- (b) When is using a BST preferred over an AVL tree?
- (c) Consider an AVL tree with n nodes and a height of h . Now, consider a single call to `get(...)`. What's the maximum possible number of nodes `get(...)` ends up visiting? The minimum possible?
- (d) **Challenge Problem:** Consider an AVL tree with n nodes and a height of h . Now, consider a single call to `insert(...)`. What's the maximum possible of nodes `insert(...)` ends up visiting? The minimum possible? Don't count the new node you create or the nodes visited during rotation(s).

6. Design

Imagine a database containing information about all trains leaving the Washington Union station on Monday. Each train is assigned a departure time, a destination, and a unique 8-digit train ID number.

What data structures you would use to solve each of the following scenarios. Depending on scenario, you may need to either (a) use multiple data structures or (b) modify the implementation of some data structure.

Justify your choice.

- (a) Suppose the schedule contains 200 trains with 52 destinations. You want to easily list out the trains by destination.
- (b) In the question above, trains were listed by destination. Now, trains with the same destination should further be sorted by departure time.

- (c) A train station wants to create a digital kiosk. The kiosk should be able to efficiently and frequently complete look-ups by train ID number so visitors can purchase tickets or track the location of a train. The kiosk should also be able to list out all the train IDs in ascending order, for visitors who do not know their train ID.

Note that the database of trains is not updated often, so the removal and additions of new trains happen infrequently (aside from when first populating your chosen DS with trains).

7. Evaluating hash functions

Consider the following scenarios.

- (a) Suppose we have a hash table with an initial capacity of 12. We resize the hash table by doubling the capacity. Suppose we insert integer keys into this table using the hash function $h(x) = 4x$.

Why is this choice of hash function and initial capacity suboptimal? How can we fix it?

- (b) Suppose we have a hash table with an initial capacity of 8 using quadratic probing. We resize the hash table by doubling the capacity.

Suppose we insert the integer keys $2^{20}, 2 \cdot 2^{20}, 3 \cdot 2^{20}, 4 \cdot 2^{20}, \dots$ using the hash function $h(x) = x$.

Describe what the runtime of the dictionary operations will over time as you keep adding these keys to the table.

8. Code analysis

For this problem, we will consider a hypothetical hash table that uses linear probing and implements the IDictionary interface. Specifically, we will focus on analyzing and testing one potential implementation of the remove method.

- (a) Come up with at least 4 different test cases to test this `remove(...)` method. For each test case, describe what the expected outcome is (assuming the method is implemented correctly).

Try and construct test cases that check if the `remove(...)` method is correctly using the key's hash code. (You may assume that you can construct custom key objects that let you customize the behavior of the `equals(...)` and `hashCode()` method.)

- (b) Now, consider the following (buggy) implementation of the `remove(...)` method. List all the bugs you can find.

```
public class LinearProbingDictionary<K, V> implements IDictionary<K, V> {
    // Field invariants:
    //
    // 1. Empty, unused slots are null
    // 2. Slots that are actually being used contain an instance of a Pair object

    private Pair<K, V>[] array;

    // ...snip...

    public V remove(K key) {
        int index = key.hashCode();

        while ((this.array[index] != null) && !this.array[index].key.equals(key)) {
            index = (index + 1) % this.array.length;
        }

        if (this.array[index] == null) {
            throw new NoSuchElementException();
        }
        V returnValue = this.array[index].value;
        this.array[index] = null;
        return returnValue;
    }
}
```

- (c) Briefly describe how you would fix these bug(s).

9. Algorithm design: easier

When writing mathematical expression, we typically write expressions in *infix* notation: in the form NUM OPERATOR NUM. An example of an expression written in infix notation is $4 + 6 * 5$. This expression evaluates to 34.

An alternative way we can write this expression is using *post-fix* notation: in the form NUM NUM OPERATOR. For example, consider the following expression written in post-fix notation:

4, 6, 5, *, +

This expression is interpreted in the following way:

- Read and store 4
- Read and store 6
- Read and store 5
- Multiply the last two stored values (and remove them from storage), then store the result
- Add the last two stored values (and remove them from storage), then store the result

The last result stored is the final “output”. In this case, the expression above also evaluates to 34.

- (a) Explain how you might apply or adapt the ADTs and data structures you’ve studied so far to evaluate an expression written in post-fix notation. Assume you accept the expression you need to evaluate as a string.
- (b) Give pseudocode for this algorithm.

10. Algorithm design: harder

- (a) Given a BST, describe how you could convert it into an AVL tree. What is the runtime of your algorithm?
- (b) Give pseudocode for an algorithm that verifies that a tree satisfies all of the AVL invariants in $\mathcal{O}(n)$ time. Assume every node object has five fields: key, value, height, left, and right.

Be sure to verify that:

- The tree is actually binary search tree
- The height information of every node is correct
- Every node is balanced

Hint: rather than trying to check all three of these things in a single pass, try writing three separate methods: one per each invariant. While it's possible to check everything in one pass, doing so will be more challenging to implement.

(As a reminder: If each method takes $\mathcal{O}(n)$ time, running all three of them will still take $\mathcal{O}(n)$ time.)

Challenge Problems

11. Random Hash Functions

In class we talked about various strategies to minimize collisions. In this question we discuss how to use randomness to “spread out” collisions from a small set of very bad inputs into a larger set of almost-always-fine inputs. The last two parts of this problem are beyond the scope of this course, but are interesting nonetheless.

For simplicity, assume our keyspace (the set of possible keys) is the set $\{0, 1, 2, \dots, 2^{30} - 1\}$. Suppose we have a hashtable of size 2^{10} . Let a be an odd integer less than 2^{30} .

Consider the hash function

$$h_a(x) = \left\lfloor \frac{(ax) \bmod 2^{30}}{2^{20}} \right\rfloor$$

Notice that the function changes depending on the value of a we choose, so this is really a set of possible functions.

- (a) Show that for any a , h_a outputs an integer between 0 and $2^{10} - 1$ (i.e. we can use this as a hash function for our table size)
- (b) Choose $a = 1$, i.e. the hash function simplifies to

$$h_1(x) = \left\lfloor \frac{x \bmod 2^{30}}{2^{20}} \right\rfloor$$

For this function, find a large set of elements that all hash to 0.

- (c) Let x, y be any of the two elements you found in the last part. Choose a few thousand values of a , and check whether $h_a(x) = h_a(y)$ for each of them (write code for this part). For what fraction of these hash functions do x, y collide? If the values of the hash function were totally random, how often would you expect collisions?

- (d) The following statement is true (explaining why is beyond the scope of the course): For any x, y if you choose a at random, the probability that $h_a(x) = h_a(y)$ is at most $2/2^{10}$.

Use this fact, or your observations in the last part, to explain why you might decide to choose a random a instead of just choosing $a = 1$ (hint: imagine you know someone is using the hash function with $a = 1$, how can you use the first part to slow their code down? Can you do the same for a random a ?)