B-Trees, LLRBs, Hashing

Discussion 07



Announcements

Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
			3/6 Project 2A Due		3/8 Lab 7 Due	
					3/15 Lab 8 Due	



Content Review



B-Trees

B-Trees are trees that serve a similar function to binary trees while <u>ensuring</u> a bushy structure (check: why don't BSTs/binary trees generally?). In this class, we'll often use B-Tree interchangeably with 2-3 Trees.



Each node can have up to 2 items and 3 children. There are variations where these values are higher, known as 2-3-4 trees (nodes can have up to 3 items and 4 children).

All leaves are the same distance from the root, which makes getting take $\Theta(\log N)$ time.



B-Trees

When adding to a B-Tree, you first start by adding to a leaf node, and then pushing the excess items (typically the middle element) up the tree until it follows the rules (max 2 elements per node, max 3 children per node).



Left Leaning Red Black Trees

LLRBs are a representation of B-trees that we use because it is easier to work with in code. In an LLRB, each multi-node in a 2-3 tree is represented using a red connection on the left side.



LLRB Rules

Each 2-3 tree corresponds to a (unique) LLRB*. This implies that:

- 1. The LLRB must have the same number of black links in all paths from root to null (not root to leaf!)
- 2. A node may not have two red children
- 3. All red links should be left-leaning
- 4. Height cannot be more than ~2x the height of its corresponding 2-3 tree
- 5. Additionally, we insert elements as leaves with red links to their parent node

All these invariants mean that sometimes our LLRB becomes unbalanced (ie. it violates a rule), so we need some way to fix that.

*2-3-4 trees correspond more generally to regular Red Black Trees, but our focus in 61B is on LLRBs.



CS61B Spring

Why root to null?

Consider the left image below: each leaf is 1 black link away from the root, but it's not a valid LLRB!

This is because when we convert it to a B-tree, the [3, 7] node will only have 2 children, not 3!

If we check for root-to-null: the right of 3 is 0 black link away, but the other nulls are 1 black link away.



LLRB Balancing Operations



Hashing

Hash functions are functions that represent objects as integers so we can efficiently use data structures like HashSet and HashMap for fast operations (ie. get, put/add).

Once we have a hash for our object, we use modulo to find out which "bucket" it goes into. For example, we can create a hash function for the Dog class by overriding Object's hashCode():

```
@Override
public int hashCode() {
    return 37 * this.size + 42;
}
```

Then, when we try to put the dog into a HashSet, the HashSet code might look something like this:

```
int targetBucket = dog.hashCode() % numBuckets;
addToTargetBucket(dog, targetBucket);
```









Hashing

In each bucket, we deal with having lots of items by chaining the items and using .equals to find what we are looking for. In a HashMap, we're specifically concerned with equality of keys in key-value pairs (in HashSet, we only have a value to compare to).

<u>**Therefore, it is important that your .equals()</u> function matches the result of comparing hashcodes - if two items are equal, they must also have the same hashcode^{**}



Hashing

The load factor tells us when we should resize. We calculate it by dividing the total number of elements added by the number of buckets we currently have. When resizing up, if the load factor exceeds some threshold, we increase the number of buckets we use in the data structure.

Because all elements were initially placed into buckets based on how many buckets were previously available, we also need to rehash all elements into a potentially new destination bucket when resizing, or else subsequent calls to get() may fail.*

* Resizing sounds like a linear-time operation...how does that affect the runtimes of our operations?



Valid vs. Good Hashcodes

Properties of a valid hashcode:

- 1) Must be an integer
- 2) The hashcode for the same object should always be the same
- 3) If two objects are "equal", they have the same hashcode
 - Check! What about the reverse?

Properties of a good hashcode:

- 1) Distributes elements evenly
 - What does this even mean?



Worksheet



inserting 18, 38, 12, 13, and 20.





inserting 18, 38, 12, 13, and 20.





Inserting 18

inserting 18, 38, 12, 13, and 20.





Inserting 38 (part 1)

inserting 18, 38, 12, 13, and 20.





Inserting 38 (part 2)

inserting 18, 38, 12, 13, and 20.





Inserting 12

inserting 18, 38, 12, 13, and 20.



CS61B Spring 2024

Inserting 13 (part 1)

inserting 18, 38, 12, 13, and 20.





Inserting 13 (part 2)

inserting 18, 38, 12, 13, and 20.





Inserting 13 (part 3)

inserting 18, 38, 12, 13, and 20.





Inserting 20.

1B 2-3 Trees and LLRBs Convert the resulting 2-3 tree to a red-black tree.





1B 2-3 Trees and LLRBs Convert the resulting 2-3 tree to a red-black tree.





1B 2-3 Trees and LLRBs Convert the resulting 2-3 tree to a red-black tree.



CS61B Spring 2024

1C 2-3 Trees and LLRBs

If a 2-3 tree has depth H (that is, the leaves are at distance H from the root), what is the maximum number of comparisons done in the corresponding red-black tree to find whether a certain key is present in the tree?



1C 2-3 Trees and LLRBs

If a 2-3 tree has depth H (that is, the leaves are at distance H from the root), what is the maximum number of comparisons done in the corresponding red-black tree to find whether a certain key is present in the tree?

2H + 2 comparisons (longest path from leaf to root * 2 items per node)













the final tree after inserting 9 to an LLRB.



After rotateRight(10)













```
public int hashCode() {
     return -1;
}
public int hashCode() {
     return intValue() * intValue();
}
// Object's hashCode() based on memory location
public int hashCode() {
     return super.hashCode();
}
// return current time as an int
public int hashCode() {
     return (int) new Date().getTime();
3
public int hashCode() {
     return intValue() + 3;
```

}



// Valid, but collisions for everything

```
public int hashCode() {
     return -1;
}
public int hashCode() {
     return intValue() * intValue();
}
// Object's hashCode() based on memory location
public int hashCode() {
     return super.hashCode();
}
// return current time as an int
public int hashCode() {
     return (int) new Date().getTime();
3
public int hashCode() {
     return intValue() + 3;
```

}



```
public int hashCode() {
    return -1;
}
public int hashCode() {
    return intValue() * intValue();
}
```

```
// Object's hashCode() based on memory location
public int hashCode() {
    return super.hashCode();
}
```

```
// return current time as an int
public int hashCode() {
    return (int) new Date().getTime();
}
```

```
public int hashCode() {
    return intValue() + 3;
}
```

// Valid, but collisions: consider -3 and 3 $\,$

// Valid, but collisions for everything



```
public int hashCode() {
    return -1;
}
```

```
public int hashCode() {
    return intValue() * intValue();
}
```

```
// Object's hashCode() based on memory location
public int hashCode() {
    return super.hashCode();
}
```

```
// return current time as an int
public int hashCode() {
    return (int) new Date().getTime();
}
```

```
public int hashCode() {
    return intValue() + 3;
}
```

// Valid, but collisions: consider -3 and 3 $\,$

// Valid, but collisions for everything

// Invalid: memory address unique per Integer



2A Hashing Are the following hashCodes valid? If they are, what are the advantages or disadvantages? public int hashCode() { // Valid, but collisions for everything

```
public int hashCode() {
    return -1;
}
```

```
public int hashCode() {
    return intValue() * intValue();
}
```

```
// Object's hashCode() based on memory location
public int hashCode() {
    return super.hashCode();
}
```

```
// return current time as an int
public int hashCode() {
    return (int) new Date().getTime();
}
```

```
public int hashCode() {
    return intValue() + 3;
}
```

// Valid, but collisions: consider -3 and 3

// Invalid: memory address unique per Integer

// Invalid: non consistent for same object



2A Hashing Are the following hashCodes valid? If they are, what are the advantages or disadvantages? public int hashCode() { // Valid, but collisions for everything

```
public int hashCode() {
    return -1;
}
```

```
public int hashCode() {
    return intValue() * intValue();
}
```

```
// Object's hashCode() based on memory location
public int hashCode() {
    return super.hashCode();
```

```
// return current time as an int
public int hashCode() {
    return (int) new Date().getTime();
}
```

```
public int hashCode() {
    return intValue() + 3;
}
```

3

// Valid, but collisions: consider -3 and 3 $\,$

// Invalid: memory address unique per Integer

// Invalid: non consistent for same object





2B Hashing

1. When you modify a key that has been inserted into a HashMap will you be able to retrieve that entry again? Explain.

2. When you modify a value that has been inserted into a HashMap will you be able to retrieve that entry again? Explain.



2B Hashing

1. When you modify a key that has been inserted into a HashMap will you be able to retrieve that entry again? Explain.

Sometimes: If the hashCode() for the key happens to change as a result of the modification (which is very likely but not guaranteed), then we won't be able to retrieve the entry in our hashtable (unless we were to recompute which bucket the new key would belong to)

2. When you modify a value that has been inserted into a HashMap will you be able to retrieve that entry again? Explain.



2B Hashing

1. When you modify a key that has been inserted into a HashMap will you be able to retrieve that entry again? Explain.

Sometimes: If the hashCode() for the key happens to change as a result of the modification (which is very likely but not guaranteed), then we won't be able to retrieve the entry in our hashtable (unless we were to recompute which bucket the new key would belong to)

2. When you modify a value that has been inserted into a HashMap will you be able to retrieve that entry again? Explain.

Always: The bucket index for an entry in a HashMap is decided by the key, not the value.



```
HashMap<String, Integer> hm = new HashMap<>();
hm.put("Hashbrowns", 7);
hm.put("Dim sum", 10);
hm.put("Escargot", 5);
hm.put("Brown bananas", 1);
hm.put("Burritos", 2);
hm.put("Buffalo wings", 8);
hm.put("Banh mi", 9);
hm.put("Burritos", 10);
```





N = 1M = 4N/M = 0.25





N = 2M = 4N/M = 0.5





N = 3M = 4N/M = 0.75









N/M = 0.375

























3B A Side of Hash Browns

Do you see a potential problem here with the behavior of our HashMap? How could we solve this?



3B A Side of Hash Browns

Do you see a potential problem here with the behavior of our HashMap? How could we solve this?

Inserting many words with the same first letter results in that letter's bucket growing very large, and our current hashing scheme means that resizing doesn't help us re-distribute the elements. We could solve this by having a better hash function!

