

Sets / Maps

Python had dictionaries: what are they useful for
Brainstorm

Map<K, V> {

boolean put (K key, V value)

boolean get (Object key) \leftarrow will return false if types don't match

boolean isEmpty()

boolean remove (Object key)

Set<K> keySet() \leftarrow opt

}

The country linked by from exam was in fact a Map

Set like a map, but has no value

~~How can we implement it?~~

Set and Map can be implemented the same way, Set has less information.

Could use an ArrayList for both.

word count

factorial

Maps

Map $\langle K, V \rangle$

~~map~~ ~~put~~

✓ put(K, V) \leftarrow returns old value
✓ get(K)

Ex: Telephone Book or Contact List
Cache

name \rightarrow phone #
input \rightarrow output

Unordered List

put $O(1)$ (at end)
get $O(n)$ sequential search

Sorted Arraylist

put $O(n)$ need to shift
get $O(\log n)$ binary search

Linked

Sorted List

put $O(n)$ need to find
get $O(n)$ can't binary search a linked list

Hashing

Idea: Use an array, but rather than sorting, the index for key k is given by a function $h(k)$.
which maps keys to ints

~~h(k)~~ is

h is the hash function
 $h(k)$ is the hash code

Ex: Map names to offices

Sam \rightarrow 318

Dave \rightarrow 312

Convert each name to an index into an array of size 5?

Brainstorm

Idea: Sum what if we have Aden, Dean and Deng

First letter Dave
David

$h(k) \% \text{table size}$

Problem: We end up with collisions ← more on this Friday

Goal: Find a good hash function

~~Let's look at~~ hash ^{codes} ~~vectors~~ should be unique (even mod table size)
hash codes should be easy to compute

Strings

$$h(s) = h(s_1 s_2 s_3 \dots s_n) = s_1 + p s_2 + p^2 s_3 + \dots + p^{n-1} s_n$$

Java uses 31

Let's look at it

Define our hash code

Other topics

Passwords

or

Salting
Verifying against the hash result

Proof of work

Spam Prevention
Bitcoin

$K \times$

$h(K+x) < n$ (i.e. begins with certain number of zeros)

Bloom Filter

~~Probabilistic Data~~
Set

might contain

Recap

$h(k)$: a function from keys to integers

Strategy when

Strategy

HashMap

hash function

1 Convert keys to integers using $h(k)$

2 Store key (and value) at location given by $h(k)$ in an array, if that location is empty

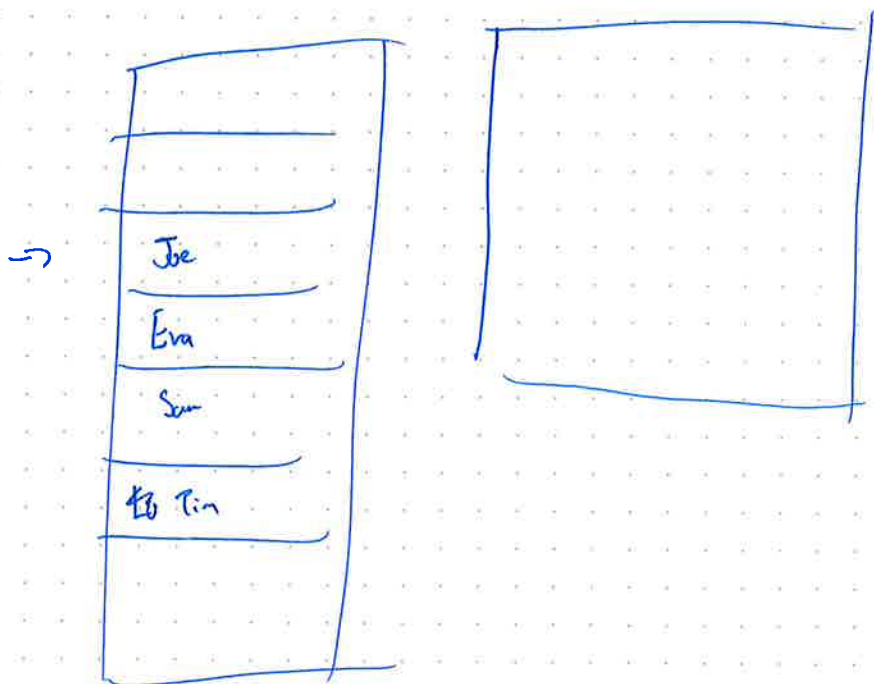
Entry(k, v) [] table;

put(k, v)

1 ~~Put~~

$i = h(k)$

if ~~store~~ (table[i] is empty or table[i].key = k)
table[i] = ~~new~~ Entry(k, v);



Analysis

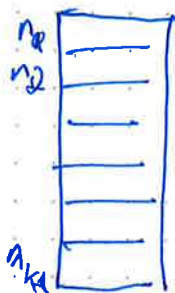
n = # of entries
 K = size of table

$$\alpha > \text{load factor} = \frac{n}{K}$$

Open addressing $\alpha \leq 1$

Chaining α can be as large as you want

Unsuccessful search



$$\frac{n_0 + n_1 + \dots + n_{K-1}}{K} = \alpha \frac{n}{K} = \alpha$$

$$\Theta\left(\frac{n}{K}\right)$$

↑ average time ↑ worst case

Successful search you go through half the chain, so ~~$\Theta(\frac{n}{K})$~~

closer to $\Theta(\frac{\alpha}{2}) \rightarrow \Theta(\alpha)$

Open addressing

Harder ^{what do you expect?} analysis

Assume we don't have clustering problems

$$P(\text{slot } h(k) \text{ is full}) = \alpha$$

Like flipping a coin with probability of heads = α

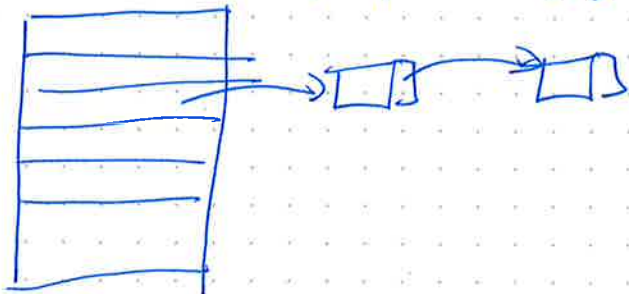
Expected ~~value~~ ^{number of steps} is $\frac{1}{1-\alpha}$

$$\Theta\left(\frac{1}{1-\alpha}\right)$$

With linear probing this looks more like $\Theta\left(\frac{1}{(1-\alpha)^2}\right)$

~~Other~~ Chaining

Other solution is to pull up a seat anyways



Problem: There might be a long chain (i.e. people sitting at your table)

Analysis

Space

Just add ^{fill in array} ~~one pointer to array~~ for open addressing
no new allocations

chaining

Constantly allocates new nodes.
all things being equal, this takes more space

Deletion

Open Addressing hard
Chaining easy

Table Size

~~No hard limit~~

Open Addressing hard limit
Chaining no hard limit
Both - Performance will deteriorate

Hashing

Collisions

Metaphor: Meeting friend for lunch in the dining hall

Strategy 1 - Open Addressing

Finding a predictable next spot to look
a. look at nearby tables

if you search nearby tables 1 by 1 linear probing $h(k) + i$

Primary
Problem 1: Clustering

Fig.

The chance of someone hitting your cluster of tables is high
so large clusters get larger more quickly

~~Leads to larger~~

b. look at tables further and further away: $h(k) + i^2$

Problem: Secondary Clustering

Two keys with the same original hash will ~~keep~~ continue to collide.

In metaphor: people with similar tastes keep showing up

c. Double hashing $h(k) + i \cdot h_2(k)$

You and your friend agree to ~~keep~~ both a location to look
and a ~~the~~ interval to schedule future searches

d. Cuckoo hashing choose two locations $h_1(k)$ and $h_2(k)$

If both locations are full, kick some one out, try
continue kicking people out.

If you reach a loop, rehash everything with new hash functions

reassign location ~~(potentially)~~
for everyone. It is ask everyone to be rehashed

Other problem - Delete

Show up and find an empty table

Solution: leave tracing

(not a problem with cuckoo hashing)

Motivation: Appointment reservation system

Single Office
Reservations for future appointments
Reserve request for appointment time t
Add t to set R if no other appointments are within 30 minutes

11:00 12:15 13:00

$|R| = n$
 $O(\log n)$ time \leftarrow target

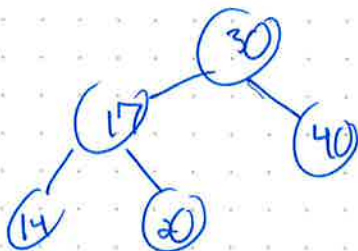
Do example

Unsorted list: Insert in $O(1)$ w/o check, but check takes $O(n)$

Sorted array: Search in $O(\log n)$ - Find smallest i such that $R[i] \geq t$
Do check in $O(1)$ - Compare $R[i]$ or $R[i-1]$ against t
Actual insertion is $O(n)$ - Need to shift

Hash ^{map} Table: can only look up exact values - cannot do range query

Binary Search Trees



Node x : Key (and values)

Pointers: parent
left
right \leftarrow These are trees!

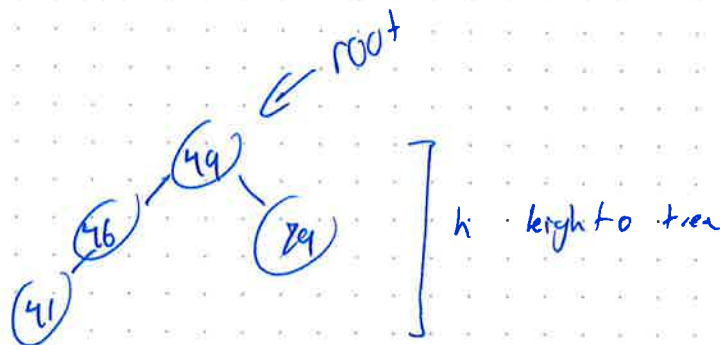
Invariant: For all nodes x , if y is in the left subtree of x then $y.key \leq x.key$

if y is in the right subtree of x then $y.key \geq x.key$

Look at example, left tree values are < 30
right tree values are > 30

Insert

49
29
46
41



$O(h)$ time

Look at code:

Traversals

Start with pre order

Use BST to sort?

Print in sorted order

find Min()

go to left until leaf $O(h)$

next Layer (t)

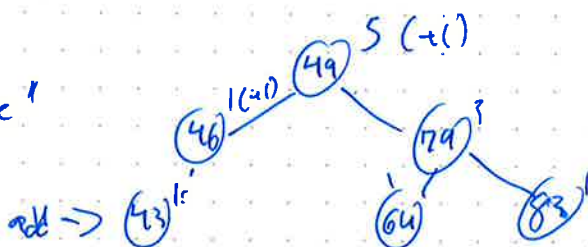
$O(h)$

Augmentation

Rank(t)

how many appointments occur before t

Store "size"



if $t < 49$
return rank to left.rank(t)
 $t = 49$
return left.size + 1
 $t > 49$
return left.size + 1 + rank(t)