Dictionaries: Hashing, **Amortization, and other magic** (there is no magic) Ilia Lebedev ilebedev@mit.edu



Hello!



Context

Part of a data structure centric Algorithms course (mid-semester) Prerequisites:

- Basic data structures
- Basic complexity

Handout: http://people.csail.mit.edu/ilebedev/facebook



Data structures cheat sheet

Insert θ(n) Array **H**(1) List O(log n) Sorted array list* $O(\log n)$ Min/Max heap $O(\log n)$ **Balanced BST**

Delete Search $\Theta(n)$ **H**(1) O(n)O(n)O(n)O(log n) O(n)O(n) $Q(\log n) O(\log n)$ **\Theta(1)**

Today: Dictionary Θ(1)

Motivating problem

Given a "document" (list of words), find the most often-occurring word.



Ready? Set. Code! (1/3) def common_word(W): counts = []best = [None, 0] $O(n^2)$ for w in W: for pair in counts: if pair[0] == w: Coooode! pair[1] = pair[1] + 1if pair[1]>best[1]: best=pair break counts.append([w, 1]) return best[0]

Ready? Set. Code! (2/3)

O(n log n)



def common word(W): SW = sorted(W)best = [None, 0]current = [None, 0]for w in SW: if w != current[0]: if current[1]>best[1]: best=current current = [w, 0]current[1] = current[1]+1 return best[0]

Ready? Set. Code! (3/3)

Use a heap to store counts? $\rightarrow O(n^2)$



Use a BST? **(augmented to store words) $\rightarrow O(n \log n)$

Can we do any better?

Finding the bottleneck

Output is a reduction of entire input $\rightarrow \Omega(n)$ complexity To find most often occurring word, we **count** all the words So far, O(n log n) at best. Given the counts, can find largest in O(n)

The dictionary ADT

• search (key) \rightarrow value

implemented

in $\Theta(1)^*$

Pre-hash functions (all keys are now ints)

 \rightarrow (pre-hash) \rightarrow Ideally: == for "same" not == otherwise

not injective!
hash("\0B") ==
hash("\0\0C")



An unbounded array "dictionary"

- all keys are
 ints ∈ Z⁺
- keys index into ∞ array A

search: $A[k] \rightarrow v$ insert: A[k] = vdelete: A[k] = None



Hashing to "compress" key space

Key space is large and sparse.

Map to dense array!

- HII

h: key space \rightarrow array indexes (mod



search: $A[h(k)] \rightarrow v$ insert: A[h(k)] = vdelete: A[h(k)] = None

Key collision

Augment a dictionary such that it maintains a list of its elements

Too

easy?

Pigeonhole principle Key space is large, but {0,1,..m-1} is small

There must exist some (many) $k_1, k_2 | h(k_1) == h(k_2)$



A hash table with chaining

Collisions will happen! Chaining: store lists (key, value)



*other solutions are very cool!

- Open addressing, Cuckoo hashing, ...

Chaining example (1/6)



Chaining example (2/6)



Chaining example (3/6)



Chaining example (4/6)



Chaining example (5/6)

Chaining example (6/6)



Time complexity with chaining

O(L_i) time to traverse the list, where L_i is |chain_i| Best case : all L_i are the same chains have N/M pairs L; should be kept short! M $\Theta(1)$ search $\rightarrow M = \Theta(N)$

The worst case

What is the worst h? O(N) worst case?!

Pick h such that this is <u>unlikely</u>. Best we can do, but it works out.

-1-1-1With SUH, the likelihood of this happening is 1/ tinv!

Many simple hash functions do well

- division method: $h(k) = k \mod m$
- multiplication method: h(k) = (k*a mod 2^w) >> (w-r),
 where w is the word width, and m = 2^r
- universal hashing: h(k) = ((a*k+b) mod p) mod m
 Demo! (random) (prime)

Various simple hash functions perform well against biased distributions. (not crypto-quality but good enough).

Taking stock

what's missing?

any solutions to challenges?



Resizing M

Recall M = $\Theta(N)$. If N = M², then L= $\Theta(\sqrt{N})$, not const! We need to resize the array as N grows!

Resizing: 1). Make a new (bigger/smaller) array 2). Make a new hash function On m 3). Re-hash everything

How long does this take?



Amortized time complexity

Resizing takes O(N) time.

- Resizing takes O(N) time. Use <u>amortization</u> keep expected (average) time low.
- $(O(1)\Theta(M) + O(N)) / \Theta(M)$ $= \Theta(M+N)/\Theta(M) = \Theta(1)$ Spread out expensive ops (amortize them)

When to resize on insert and by how much?



No, $\Theta(M)$ inserts must happen before next resize. m $\rightarrow 2m$

at least m/2 inserts before next resize

Resizing on deletions

- Tricky case: [ins, del, ins, del, ...] at threshold
- a resizing *del* must allow $\Theta(M)$ *ins* before resizing
- $m \rightarrow m/2$ when n < m/4

 $m \rightarrow m/2$ when n < m/2 is BAD!



Hash table with chaining (1/5)

Now, lets put it all together!

```
class HashTable:
 def init (self):
    self.min m = 4
    self.m = self.min m
    self.n = 0
    self.D = [ [] ] * self.m
    self.h =
     lambda k : ((k*73+37)%91)%self.m
```

Hash table with chaining (2/5)

def search (self, k): chain = self.D[self.h(k)] for (Dk, Dv) in chain: if Dk == k: return Dv return None

Hash table with chaining (3/5)

def insert (self, k, v): chain = self.D[self.h(k)] for pair in chain: if pair[0] == k: pair[1] = vreturn chain.append([k, v]) self.n = self.n + 1if self.n > self.m: self. resize(2*self.m)

```
Hash table with chaining (4/5)
def delete ( self, k ):
  chain = self.D[self.h(k)]
  m = self.m
  for pair in chain:
    if pair[0] == k:
      pair[0] = chain[-1][0]
      pair[1] = chain[-1][1]
      chain.pop()
      self.n = self.n - 1
      if (m > self.min m) and (self.n < m/4):
        self. resize(m/2)
      return
```

Hash table with chaining (5/5)

```
def resize ( self, new m ):
  elements = []
  for chain in self.D:
     for element in chian:
       elements.append( element )
  self.D = [ []*new m ]
  self.m = new m
  for (k,v) in elements:
     self.insert( k, v )
```



uch complexity

so hash

much table

any questions?

WOW

WOW

wow