

WEEK 4 L1

HASH TABLES

Purpose: immediate random access, through keys

Insert: add new record

Delete: delete existing record

Lookups: check for a particular record

Using key

ONLY THESE OPERATIONS, BUT GUARANTEED O(1) TIME (if implemented well)

- (1) Properly implemented
- (2) Non pathological

APPLICATION: DE-DUPLICATION

Input: "stream" of objects

Goal: ignore duplicates

hashable Hash lookups.

If not found, insert x into M

APPLICATION: 2-SVM PROBLEM

Input: Unsorted array of n integers. Target sum t

Goal: determine whether two numbers in Array $x + y = t$

Name solution: $O(n^2)$ ($n(n-1)$)

SUM PROG CND

- Even better:
- ① insert elements of A in H . $O(n)$
 - ② for each x , look in H for $t-x$ $O(n)$

FURTHER OBVIOUS APPLICATIONS

- historical application: symbol tables in compilers
- blocking network traffic
- search algorithms (game tree exploration), can use to avoid exploring more than once

WEEK 4 L2

HASH TABLE IMPLEMENTATION

High Level rules:

Setup: Universe U (generally, really big)

Goal: Maintain evolving set $S \subseteq U$ (of reasonable size)

Naive Solution ①: Array based solution indexed by n ($O(1)$ ops, $O(n)$ space)

Naive Solution ②: List based solution ($O(1)$ space, $O(s)$ lookups)

Solution:

- ① Pick $n = \# \text{ of buckets}$ to store data in with $n \geq |S|$
splits out position in say
- ② Choose a hash function $h: U \rightarrow \{0, 1, 2, \dots, n-1\}$
for simplicity we'll assume S doesn't vary too much
- ③ Use array A of length n , store x in $A[h(x)]$

What happens when there are collisions? $h(\text{"Bob"}) = h(\text{"Alice"})$

Think of birthdays! 23 people!

RESOLVING COLLISIONS:

Collision: $h(x) = h(y)$ with $x, y \in U$.

Solution #1: Separate chaining. Keep linked list in bucket, perform operations on lists.

Solution #2: Open Addressing (one object per bucket), h now specifies probe sequence

~~Double hashing~~: use buckets as hash lists themselves, takes lots of time!

(Linear probing, quadratic probing)

WEEK 4 C3

HASH TABLES IMPLEMENTATION DETAILS

WHAT IS A GOOD HASH FUNCTION?

Note: in hash table with chaining, insert $O(1)$ most of the time

Performance depends on the choice of hash function!

Properties of good hash function:

- ① Should spread data out well (gold standard: completely random)
- ② Should be simple enough, not take too much time (must be $O(1)$)

EXAMPLES OF BAD HASH FUNCTIONS

Ex: Keys = phone numbers $h(x) = 10^{10}$

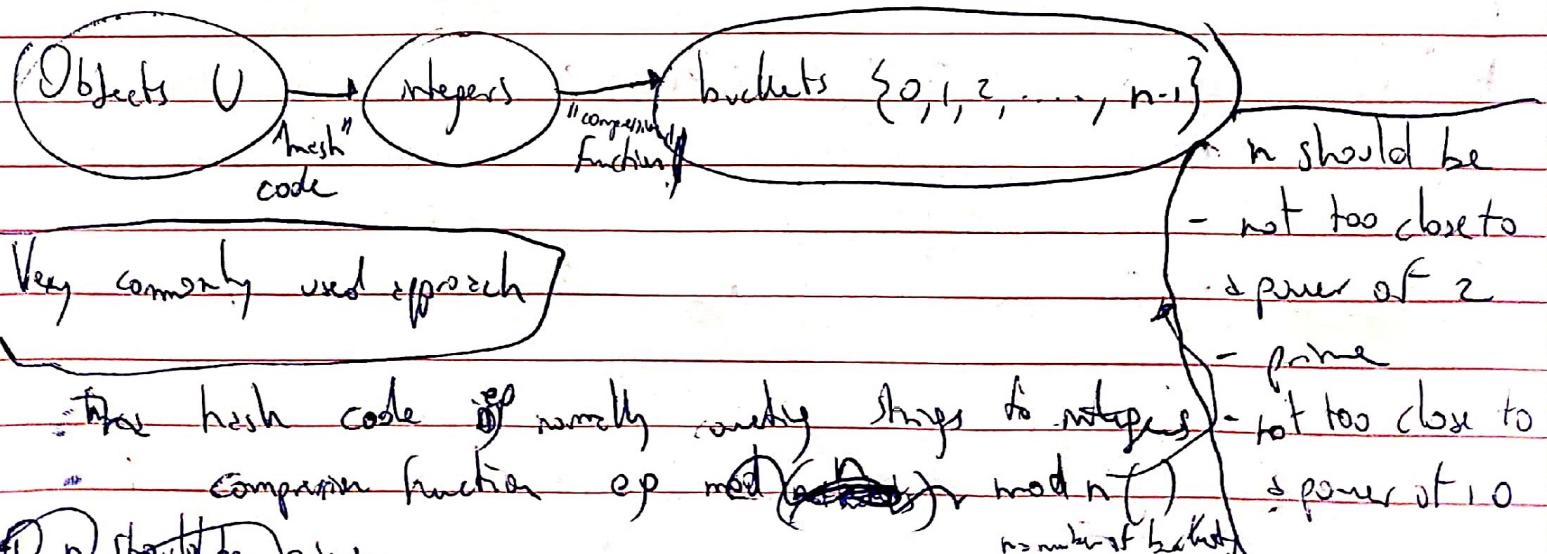
terrible: $h(x) = \text{1st 3 digits of } x \quad (n = 10^3)$

mediocre: $h(x) = \text{last 3 digits of } x \quad (n = 10^3)$

Example: keys = memory locations (power of 2)

bad hash func: $h(x) = \text{mod 1000}(x)$ ('all slot buckets are empty')

QUICK AND DIRTY HASH FUNCTIONS



WEEK 4 L4

THE LOAD OF A HASH TABLE

Definition: The load factor of a hash table is:

$$\alpha = \frac{\text{# of objects in buckets}}{\text{# of buckets}}$$

Conditions for good performance:

- ① $\alpha = O(1)$, else won't run in constant time. (NEED TO BE ACC
TO COMPUTE α in O(1))
- ② for open addressing, $\alpha < 1$ middle

For good performance, we need to control load

Also, THE PERFECT HASH FUNCTION DOES NOT EXIST. FOR EVERY $H(x)$ THERE IS A PATHOLOGICAL DATA SET

PATHOLOGICAL DATA IN THE REAL WORLD

~~Remember~~ Mem Point : Can analyze several real-world systems by exploiting badly designed hash functions.

If the hash function ~~is~~ ^{is} open source and ~~too~~ simple, ~~it's~~ easy anyone can read it and reverse engineer it.

SOLUTIONS

- ① Use a cryptographic hash function which can't be reverse engineered
- ② Use randomization: design a family H of hash functions and pick one at random

BLOOM FILTERS

A variant on hash tables, save space, have errors.

PROS

- Space

CONS

- Can't store associated object
- No deletions
- Small false positive probability

Digital APPLICATIONS

Digital: early spell checkers

(anonymize): List of forbidden passwords

Modern: Network routers

- limited memory, need to be super fast still

WHAT IS A BLOOM FILTER?

Ingredients:

- An array of n bits. $\frac{n}{|S|} = \# \text{ of bits per object in the set}$
- K hash functions h_1, h_2, \dots, h_K . $K = \text{small constant}$

for $i = 1, 2, \dots, K$, $A[h_i(x)] = 1$, whether or not it was 1 already

Lookup (x) return TRUE if $A[h_i(x)] = 1$ for every $i = 1, 2, \dots, K$

There are false positives, but no false negatives and it saves space

WEEK 4 L 9

HEURISTIC ANALYSIS OF BLOOM FILTERS

There is a trade off between % of false positives and space saved.

After inserting object s , $1 - \left(1 - \frac{1}{n}\right)^{K(S)}$ probability for a given bit to have been set to 1

$$\text{prob} \in 1 - e^{-\frac{K}{b}} \quad b = \frac{n}{|S|}, \text{ number of bits per object}$$

This means the false positive prob. is bounded by $(1 - e^{-\frac{K}{b}})^K$

how to find K ? for fixed b , the false positive probability is minimised by setting $K \approx (\ln 2) \cdot b$

Example : $b=8$, $K=5$ or 6 , error prob. $\approx 2\%$

Problem Solving

$$h(n \text{ keys}) \rightarrow \sqrt{\log_2 m}$$



like with the n tasks for m series

$\Omega(n)$

$\Omega(1)$



$$\frac{n(n-1)}{2} \cdot \frac{1}{m} = \frac{n(n-1)}{2m}$$

