

COP 3502: 11/12/21

Hash Tables

- inserts
 - deletes
 - searches
- } faster $O(\lg n)$

* Goal: Average case $O(1)$ for each
not worry too much about worst case.
(Impossible to make an input force worst case.)

Hash function

Input: Anything (arbitrary string)

Output: fixed integer in some range
 0 to $n-1$, for some given n .

Many-to-one (multiple inputs may map to
the same output.)

Properties of a good hash function

1) each output value is equally likely.

2) Given some output value it should be difficult
to construct an input that "creates" that
value. [Given y , it's hard to find x s.t.
 $f(x) = y$.]

3) If we make small changes to x , then
(say x and x'), there should many
changes btw $f(x)$ and $f(x')$.

How to create a hash table w/o hash function?

Array of size n

$$\text{Calculate } f(\text{"cat"}) = 4$$

2. Place cat in index 4.

$$f(\text{"elephant"}) = 1$$

$$f(\text{"dog"}) = 0$$

0	dog
1	elephant
2	
3	
4	cat
\vdots	
$n-1$	

Problem?

What if $f(\text{"giraffe"}) = 4$ but cat's already there?

COLLISION

4 methods of dealing w/ collisions

1) Overwrite data, lose the old data.

- ADVANTAGE (EASY)

- DISADVANTAGE (LOSE STUFF)

2) Linear Probing

If a slot, k , is full when you try to insert go to slot $(k+1) \bmod n$ and see if it's empty. Continue until you find the first empty slot.

- ADVANTAGE (DON'T LOSE STUFF)

- DISADVANTAGE (FINDING ISN'T AS FAST)

When looking for my giraffe, I can't stop looking after index 4 until I find an empty slot, so they this could slow me down a lot.

- Clustering

As there are more collisions, the probability of hitting a cluster increases and linear probing ensures those clusters will grow, if h.t.

3) When looking for new places go to these indexes:

$$\begin{aligned} & k \\ & (k+1) \text{ oloop} \\ & (k+4) \text{ oloop} \\ & (k+9) \text{ oloop} \\ & (k+16) \text{ oloop} \\ & \vdots \\ & (k+i^2) \text{ oloop} \\ & \left. \begin{array}{l} (k+(\frac{i-1}{2})^2) \text{ oloop} \end{array} \right\} p \end{aligned}$$

Quadratic Probing

- (a) table size prime
- (b) less than half full

Prove each # of this list is unique.

Assume the opposite that

$$\begin{aligned} & k+i^2 \equiv k+j^2 \pmod{p} \\ \Rightarrow & (i-j)(i+j) \equiv 0 \pmod{p} \\ & i^2 - j^2 \equiv 0 \pmod{p} \\ & (i-j)(i+j) \equiv 0 \pmod{p} \\ \Rightarrow & p \mid [(i-j)(i+j)] \end{aligned}$$

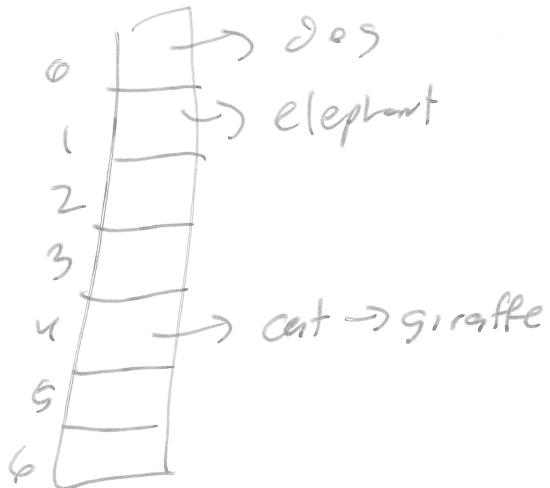
Since $p \in \text{Prime}$, $p \mid (i-j)$ or $p \mid (i+j)$

$$\begin{array}{c} \downarrow \\ \text{FALSE} \\ \text{because } i \neq j \text{ and } |i-j| < \frac{p}{2} \end{array} \quad \begin{array}{c} \downarrow \\ \text{FALSE} \\ 0 < i+j < p. \end{array}$$

For both linear + quadratic probing, double table size when it gets half full

4) Separate Chaining Hashing

At end index have a linked list



$$\begin{aligned}f(\text{cat}) &= 4 \\f(\text{dog}) &= 0 \\f(\text{elephant}) &= 1 \\f(\text{giraffe}) &= 4\end{aligned}$$

Idea is runtimes are probability based on the probability of searching in each slot. Longer ~~lists~~ lists \Rightarrow worse run times.

Why would someone use Quadratic Probing over this? (Index math in an array is faster than mallocing + too following links in lists.)

Hash Functions

$f(\text{"cat"}) = (\text{'c'} + \text{'a'} + \text{'t'})$ bad, many values map to the same spot + for most words sum of ascii values clusters in the 200-1500 range.

Improvement: $f(\text{"cat"}) = (\text{'c'} \times 128^2 + \text{'a'} \times 128^1 + \text{'t'} \times 128^0)$ on writing a string in some "base"

$$p(\text{search } L_{\text{list}}) = \frac{\# \text{ words in list}}{\text{Total words}}$$

$$\text{Avg Run Time (List)} = \frac{\text{list len} + 1}{2}$$

$$\sum_{\substack{\text{lists} \\ i=0}}^{\substack{100}} \left(\sum_{\substack{\text{freq[i]} \\ \text{lists}}} \times \frac{\# \text{ word is list}}{\text{TW}} \times \frac{(\text{list len} + 1)}{2} \right)$$

ans for each
list of
length i

TW

num lists
length i