

Lec9_Table doubling, Karp-Rabin

Size of Table

- Problems
 - want $m = \Theta(n)$ at all times
 - don't know how large n will get at creation
 - m too small \Rightarrow slow; m too big \Rightarrow wasteful
- Idea
 - Start small(constant) and grow (or shrink) as necessary
- Rehashing
 - To grow or shrink table hash function must change (m, r)
 \Rightarrow must rebuild hash table from scratch
for item in old table : \rightarrow for each slot, for item in slot
Insert into new table
 $\Rightarrow \Theta(n+m)$ time = $\Theta(n)$ if $m = \Theta(n)$
- How fast to grow?
When n reaches m
 - $m += 1$?
 \Rightarrow rebuild every step
 $\Rightarrow n$ inserts cost $\Theta(1+2+3+ \dots + n) = \Theta(n^2)$
 - $m *= 2$? $m = \Theta(n)$ still ($r += 1$)
 \Rightarrow rebuild at insertion 2^i
 $\Rightarrow n$ insert cost $\Theta(1+2+4+8+ \dots + n)$
 $\Rightarrow \Theta(n)$

Amortized Analysis

- Amortized Analysis
 - Operation has amortized cost $T(n)$ if k operations cost $\leq k * T(n)$
 - “ $T(n)$ amortized” roughly means $T(n)$ “on average”, but averaged over all ops.
 - e.g. inserting into a hash table takes $O(1)$ amortized time
- Back to hashing
 - Maintain $m = \Theta(n) \Rightarrow \alpha = \Theta(n) / n = \Theta(1) \Rightarrow$ support search in $O(1)$ expected time
- Delete
 - Also $O(1)$ expected as is
 - Space can get big with respect to n e.g. n X insert, n X delete
 - Solution : when n decrease to $m/4$ shrink to half the size $\Rightarrow O(1)$ amortized cost for both insert and delete \rightarrow analysis is h
- Resizable Array :
 - same trick solves Python “list” (array)

- \Rightarrow list.append and list.pop in $O(1)$ amortized

String Matching

Given two strings s and t , does s occur as a substring of t ? (and if so, where and how many times?)

- Simple Algorithm

```
any(s==t[i:i+len(s)] for i in range(len(t) - len(s))
→  $O(|s|)$  time for each substring comparison
⇒  $O(|s|*(|t|-|s|))$  time
=  $O(|s|*|t|)$  : potential quadratic
```

- Karp-Rabin

- Compare $h(s) == h(t[i : i + len(s)])$
- If hash values match, likely so do strings
 - can check $s == t[i : i + len(s)]$ to be sure ~ cost $O(|s|)$
 - if yes, found match -- done
 - If no, happened with probability $< 1/|s|$
 - \Rightarrow expected cost is $O(1)$ per i
- need suitable hash function
- expected time = $O(|s| + |t| * \text{cost}(h))$
 - Naively $h(x)$ costs $|x|$
 - We'll achieve $O(1)$
 - ideae : $t[i : i + len(s)] \approx t[i+1 : i + 1 + len(s)]$

Rolling Hash ADT

Maintain string x subject to

- $r()$: reasonable hash function $h(x)$ on string x
- $r.append(c)$: add letter c to end of string x
- $r.skip(c)$: remove front letter from string x , assuming it is c

Karp-Rabin Application

```
for c in s : rs.append(c)
for c in t[:len(s)] rt.append(c)
if rs() == rt()
for i in range(len(s), len(t));
    rt.skip(t[i-len(s)])
    rt.append(t[i-len(s)])
    if rs() == rt()
        check whether  $s == t[i-len(s)+1 : i+1]$ 
        if equal : found match
        else : happens with probability  $\leq 1/|s|$ 
```

- Running time : $O(|s|+|t| + \text{\#match} \cdot |s|)$

Data Structure

Treat string x as a multi digit number in base a where a denotes the alphabet size, e.g, 256

- $r() = u \bmod p$ for prime $p \sim |s|$ or $|t|$ (division method)
- r stores $u \bmod p$ and $|x|$ (really $a^{|x|}$) not u
 \Rightarrow smaller and faster to work with ($u \bmod p$ fits in one machine word)
- $r.append(c)$

$$: (u \cdot a + \text{ord}(c)) \bmod p = [(u \bmod p) \cdot a + \text{ord}(c)] \bmod p$$

- $r.skip(c)$

$$[u - \text{ord}(c) \cdot (a^{|u|-1} \bmod p)] \bmod p$$

$$: [(u \bmod p) - \text{ord}(c) \cdot (a^{|x|-1} \bmod p)] \bmod p$$