

# Lecture 5: Sets & Dictionaries via... hashing (later) ①

Abstract Data Type: Set of strings [say]

Ops: (Initialize)

- Insert(s) // s a string
- IsMember(s) // true/false  
Lookup

[e.g. passwords, names, tweets; many possible not all known in advance]

[Do example, human alg., requesting

~~from~~ from class TV/movie names show

Write on slips.

Store. First unordered. Then

ordered....

Linked-list?

BST...?

Could get into those, ...?

Binary Search Tree:  $O(\log n)$  time ops if  $n$  elts.

[kind of complicated.

Can we exploit  $O(1)$ -time array lookup?]

• If you know all elts  $s$  are 3-digit #'s...

alloc a Boolean array of len. 1000!

But what if they're 10-digit nums: 10B-size array :-

280-char tweets ☹

Idea: Hashing: invent "hash" fcn:  $h: \{\text{all strings}\} \rightarrow \{0, 1, \dots, n-1\}$

• Initialize:  $T[0 \dots n-1] := \text{false}$

• Insert(s):  $T[h(s)] := \text{true}$

• IsMember(s): return  $T[h(s)]$ .  
Lookup

$n =$   
↑  
out of space you can afford.

Issues!



Issue 1: Where/how to get h?

Issue 2: Collisions: #strings = ∞. Even if len. bound, probably #strings >> n.

∴ ~~same~~ (PHP) multiple strings  $s_1, s_2$  have  $h(s_1) = h(s_2) \Rightarrow$  erroneous lookups, potentially

Solution 2: "Chaining" [other possibls...]

$T[\cdot]$  stores linked lists, not just true/false.

Init  $\rightarrow T[i] := \text{null}$

Insert(s)  $\rightarrow$  insert 's' at end of list  $T[h(s)]$

IsMember: look thru list at  $T[h(s)]$  for 's'  
Lookup(s)

Bonus: can store any data/val v together with "key" s, get a "dictionary" data type, like python's "MyDict[s] = v".

~~Worst-case time for insert/lookup~~ Space (mem.) usage:  $O(n \leftarrow \text{len. of } T[\cdot], \text{ you choose})$   
 $+ M \leftarrow \# \text{ items inserted}$   
 $(\cdot \leftarrow \text{size of item, but say } O(1))$

Worst-case time for insert/lookup(s):  $O(\text{length of } T[h(s)])$

maybe m  $\therefore$  [BST has  $O(\log m)$ ]

$\rightarrow$  could set  $n=1, h(s)=0 \forall s \rightarrow$  just a linked list!



Idea: Suppose  $h(s)$  was somehow a "random-looking" # from  $0 \dots n-1$  (but ~~was~~ actually deterministic) ③  
 $\hookrightarrow$  e.g. (ascii int. represent of  $s$ ) mod  $n$  ???

Hope: after inserting  $s_1, \dots, s_m$ ,  
the #'s  $h(s_1), \dots, h(s_m)$  are "random-ish",  
so "average" len. of a list is  $\frac{m}{n}$ . [is that random?]

$\Rightarrow$  "average-case" insert/lookup =  $O(1 + \frac{m}{n})$  [curr # elts]

$\Rightarrow$  If you estimate set will have  $\approx m$  elts, ?  
can choose  $n = m$ , have space  $O(m)$  😊  
"avg-case time"  $O(1)$  😊



[This is the key idea. But many q's to answer...]

- $h(s) = \text{int}(s) \bmod n$  is bad... // strings ending w/ same few chars map to same thing...
- how to make a "random-ish"  $h$ ?
- what is "avg-case", formally? any worst-case guarantees?
- what if you can't est.  $m$  in advance — or,  
what if you know elts, in advance, just want to do lookups?
- methods other than chaining to handle collisions?



Random-ish hash fun?

[let's think on this first!]

① Dream scenario:

Some oracle does...

for  $s$  in all strings (!!!)  
 $h(s) := \text{Rand.Int}(0 \dots n-1)$

And then you can "magically" use  $h$  in  $O(1)$  time,  
no ~~space~~ space.

② Reality attempt, python:

$h = \text{ord}(s[0]) \ll 7$   
for char in s

$h = \text{E\_mul}(1000003, h) \wedge \text{ord}(\text{char})$  // maybe earlier python, now "Siphash"  
 $h = (h \wedge \text{len}(s)) \% n$

// this is completely deterministic, but "seems sorta random".

[But knowing python does this, an evil person could keep inserting  $s$ 's that hash to same #...]

"hash flooding" or "hash DoS"

changed in 2012  
"The Power of Evil Choices in Bloom Filters."

~~Principled way~~

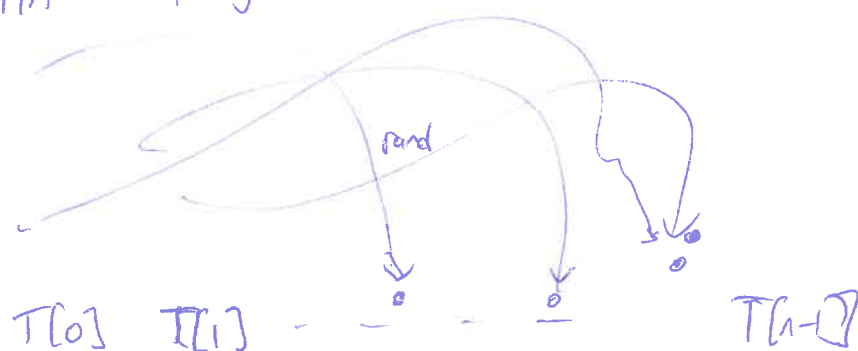
Today: "SUHA": Simple Uniform Hashing Assumption:  
do ②, pretend it's ①

[Good for analyzing hash tables in ideal circs.]

Later: Principled ways to provably achieve things achievable with SUHA.

Under SUHA: if you insert  $m$  items into size- $n$  table...

$h(s_1)$   
 $h(s_2)$   
 $h(s_3)$   
 $h(s_4)$



"Balls and Bins"  
[a very common scenario]



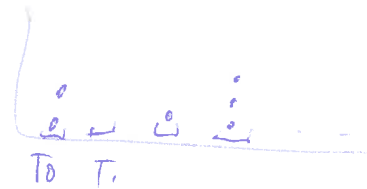


"Average load / length":  $\lambda = \frac{m}{n}$ .

(5)

Say you now do an unsuccessful lookup(s)

$\therefore s$  not in table, already  $h(s) \sim \text{RandInt}(0 \dots n-1)$



$\therefore E[\text{length of } T[h(s)]]$   
 $= \text{avg \{ bin sizes \}} = \lambda \rightarrow \text{time } O(\lambda)$

Say "successful lookup"  $s$  in table. Say  $h(s) = i$ .

All  $m-1$  other items act like rand-balls, so

$E[\text{length of } T[h(s)]] = E[\# \text{ balls in bin } i \text{ after } m-1 \text{ throws}]$   
 $= \frac{m-1}{n} \leq \frac{m}{n} = \lambda \rightarrow$

$\therefore$  space =  $O(n)$ , time =  $O(1 + \lambda) = O(1 + \frac{m}{n})$  or average  $\int \leftarrow$  rigorously def'd under-suit!

$m=n$ : popular choice. Space:  $O(m)$ , or time  $O(1)$  on avg

Say  $m=n$ .

Q: What is worst-case lookup time ... on average?

load of most loaded bin.

$\lambda=1$  but ... very unlikely all balls in diff bins!!!

$\text{Pr}[\text{first 2 balls go into } T[0]] = \frac{1}{n^2}$

$\text{Pr}[\text{exactly 2 balls go into } T[0]] = \frac{\binom{n}{2} \cdot \left(\frac{1}{n}\right)^2 \cdot \left(1 - \frac{1}{n}\right)^{n-2}}{\underbrace{\frac{n(n-1)}{2} \leq \frac{n^2}{2} \leq 1}} \leq \frac{1}{2}$

(???) [What is this if  $n=1000$ ]

$\text{Pr}[\text{exactly 3 go into } T[0]]$

$= \binom{n}{3} \left(\frac{1}{n}\right)^3 \left(1 - \frac{1}{n}\right)^{n-3} \leq \frac{n(n-1)(n-2)}{3!} \cdot \frac{1}{n^3} \cdot 1 \leq \frac{1}{3!}$



Pr[exactly 10 go into  $T[0]$ ]  $\leq \frac{1}{10!}$

(6)

Pr[ $\geq 10$  go into  $T[0]$ ]

$= \Pr[\{=10 \text{ go in}\} \cup \{=11 \text{ go in}\} \cup \{=12 \text{ go in}\} \cup \dots]$   $\leq \Pr[=10 \text{ in}] + \Pr[=11 \text{ in}] + \dots$  "Union Bound"

$\leq \frac{1}{10!} + \frac{1}{11!} + \frac{1}{12!} + \frac{1}{13!} + \dots$   
 $= \frac{1}{10!} (1 + \frac{1}{11} + \frac{1}{11 \cdot 12} + \frac{1}{11 \cdot 12 \cdot 13} + \dots)$   
 $\leq \frac{1}{10!} (1 + \frac{1}{11} + \frac{1}{11^2} + \frac{1}{11^3} + \dots)$   
 $= \frac{1}{10!} (\frac{1}{1 - \frac{1}{11}}) = \frac{1}{10!} \cdot \frac{11}{10}$

Pr[ $\geq k$  go into  $T[0]$ ]  $\leq \frac{1}{k!} \cdot \frac{k+1}{k} \leq \frac{2}{k!}$   $\leq \frac{2}{2^{k-1}} = \frac{4}{2^k}$  (lazy/stupid!)

Pr[ $\geq k$  in  $T[0]$ ]  $\cup$  [ $\geq k$  in  $T[1]$ ]  $\cup \dots \cup$  [ $\geq k$  in  $T[n-1]$ ]

$\leq \Pr[\geq k \text{ in } T[0]] + \dots + \Pr[\geq k \text{ in } T[n-1]]$  U.B.

$\leq \frac{2}{k!} + \frac{2}{k!} + \dots + \frac{2}{k!} \leq \frac{4}{2^k} n \rightarrow \leq \frac{1}{2^{20}} = \text{million}$   
 if  $n \leq 2^k$

How much is  $k!$ ?

$100! = 100 \cdot 99 \cdot 98 \cdot 97 \dots \cdot 2 \cdot 1$

$\geq 2 \cdot 2 \cdot 2 \dots \cdot 2 = 2^{99}$

or  $\geq 3 \cdot 3 \cdot 3 \dots \cdot 3 = 3^{98}$   $\leftarrow$  which is bigger? [way bigger]

or  $\geq 4^{97}$  or  $5^{96}$  or  $\dots \geq 50^{51} \geq 50^{50}$

(Not optimal, should do  $\frac{100}{e}$ , not  $\frac{100}{2}$ )

$\log(3^{98}) = 98 \cdot \log 3 = 98 \cdot (1.6)$   
 $\log(2^{99}) = 99 \log 2 = 99$

So  $k! \geq (\frac{k}{2})^{k/2}$   $\log(k!) \geq \frac{k}{2} \log(\frac{k}{2}) = \Omega(k \log k)$

$2^{22} \cdot n \leq 2^{\Omega(k \log k)}$

$\Rightarrow k \leq \frac{\log n}{\log \log n}$

max load  $\leq \log n + 22$  except w.p.  $1/\text{million}$

Q: how strongly is this achieved?

