

# Lecture 6: Max Load, Power of 2 choices, Bloom Filters.

(1)

[When we last left our heroes...]

Analyzing typical worst-case insert/lookup time (under SUHA)  
when  $m=n$  items hashed to  $n$  slots.

typical max load  $\approx$  when  $m=n$  balls thrown into  $n$  bins.



Recall: [we ended with]  $\Pr[T[0] \text{ gets exactly } k \text{ balls}] \leq \frac{1}{k!}$

$\Pr[T[0] \text{ gets } \geq 10 \text{ balls}] \leq \Pr[10] + \Pr[11] + \dots$  "union bound"  
 $\leq \frac{1}{10!} + \frac{1}{11!} + \frac{1}{12!} + \dots$   
 $= \frac{1}{10!} \left( 1 + \frac{1}{11} + \frac{1}{12 \cdot 11} + \frac{1}{13 \cdot 12 \cdot 11} + \dots \right)$   
 $\leq \frac{1}{10!} \left( 1 + \frac{1}{11} + \frac{1}{11^2} + \frac{1}{11^3} + \dots \right)$   
 $= \frac{1}{10!} \left( \frac{1}{1-1/11} \right) \leftarrow \frac{11}{10}$   
 $= \frac{11}{10} \cdot \frac{1}{10!}$  [just slightly bigger than  $\frac{1}{10!}$ ]

$\Pr[T[0] \text{ gets } \geq k \text{ balls}] \leq \frac{k+1}{k} \cdot \frac{1}{k!} \leq 2 \cdot \frac{1}{k!}$  [more like  $\approx 1 \cdot \frac{1}{k!}$ , but oh well!]

$\Pr[T[1] \text{ gets } \geq k \text{ balls}] \leq \frac{2 \cdot 1}{k!}$

$T[2] \dots$

$\Pr[|T[0]| \geq k \text{ OR } |T[1]| \geq k \text{ OR } \dots \text{ OR } |T[n-1]| \geq k] \leq$

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

$\leq \frac{2}{k!} + \frac{2}{k!} + \dots + \frac{2}{k!} = \frac{2}{k!} \cdot n$

union bound again

$\Pr[\text{max load} \geq k]$

Say  $n = 1000$ .

$$\Pr[\text{max load} \geq k] \leq \frac{2000}{k!}$$

(2)

$$\geq 6] \leq \frac{2000}{720} \approx 2.78$$

$$\geq 7] \leq \frac{2000}{5040} \approx 0.4$$

$$\geq 8] \leq \frac{2000}{40320} \approx 0.05$$

Say  $n = 10^6$ , want to say " $\Pr[\text{max load} \geq k] \leq \frac{1}{500}$ " (say)

What  $k$ ? Need  $\frac{2 \cdot 10^6}{k!} \leq \frac{1}{500} \Leftrightarrow k! \geq 10^9 \Leftrightarrow k \geq 12$ .

"For 1M balls  $\rightarrow$  1M bins, max load  $\leq 12$  except w.p.  $\leq \frac{1}{500}$ ."

For general  $n$ , need  $k! \geq 1000n$ . [Inverse factorial?]

~~Fact: 22! has 22 digits. 23! has 23 digits. 24! has 24 digits.~~

How big is  $k!$  ??  ~~$\log(k!) \approx \log 1000 + \log n$~~  Skip

~~$\log(k!) \approx \log k + \log(k-1) + \dots + \log 1$~~

# binary digits  $\sim \log_2(k!) = \lg(k \cdot (k-1) \cdot (k-2) \dots 3 \cdot 2 \cdot 1)$

$$= \lg k + \lg(k-1) + \lg(k-2) + \dots + \lg\left(\frac{k}{2}\right) + \dots + \lg 1$$

$$\geq \lg \frac{k}{2} + \lg \frac{k}{2} + \lg \frac{k}{2} + \dots + \lg \frac{k}{2}$$

$$= \frac{k}{2} \cdot \lg \frac{k}{2} = \frac{k}{2} (\lg k - 1) \geq \frac{k}{2} \cdot \frac{\lg k}{2}$$

if  $\lg k \geq 2$   
 $\Leftrightarrow k \geq 4$

$$= \frac{1}{4} k \lg k$$

$$= \Omega(k \lg k)$$

" $k!$  is a  $\Theta(k \log k)$ -digit #"

Fact:  $22!$  has 22 digits.  $\therefore 23! \geq 10 \cdot 22!$  has  $\geq 23$  digits  
 $24! \geq 10 \cdot 23!$  has  $\geq 24$  digits.

$k!$  has  $\geq k$  digits (ex: more like  $\Theta(k \log k)$  digits)

$n$  has  $\sim \log_{10} n$  digits. So if  $k \geq \log_{10} n + 3$ ,  $k! \geq 1000n$ .  
 $\therefore$  "For  $n$  balls  $\rightarrow n$  bins, max load  $\leq \Theta(\log n)$  except w.p.  $\leq \frac{1}{1000}$ ." (ex:  $\leq O\left(\frac{\log n}{\log \log n}\right)$ )

# "The power of 2 choices"

[still with "SHTA"]

③

Idea: Get hold of two ("independent") hash functions,  $h_1, h_2$ :

- Insert( $s$ ): Look at  $T[h_1(s)], T[h_2(s)]$ . linked lists  $\{strings\} \rightarrow \{0, 1, \dots, n-1\}$   
Append  $s$  to whichever is shorter. [break ties arbit.]
- Lookup( $s$ ): Look for  $s$  in both  $T[h_1(s)], T[h_2(s)]$ .

Can't really make things more than 2x worse.  
Surprise: ~~max~~ worst-case time goes way down, typically

Balls & bins ver: To "throw" a ball, pick 2 bins at random, put ball in less-loaded bin.

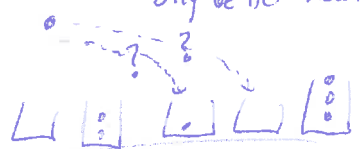
Theorem: For  $m=n$  balls, with ~~high~~ high probability, max load is  $O(\log \log n)$ . [Way better than  $O(\log n)$ . "exponentially"! For every  $\#n$  in universe,  $\log \log n \approx 8$ ]

[Proof is a little elaborate, so I'll just give...]

Idea of why: After throwing  $n$  balls, let  $\alpha_2 =$  fraction of bins with  $\geq 2$  balls

Claim:  $\alpha_2 \leq \frac{1}{2}$ . Because if  $\geq \frac{1}{2}$  bins have  $\geq 2$  balls,  $\Rightarrow$   $> \frac{n}{2}$  balls,  $\Rightarrow$   $> n$  balls,  $\Rightarrow$  ~~impossible~~.

Say we've thrown some of the balls  $\alpha_2 \leq \frac{1}{2}$  now. [It's even  $\leq \frac{1}{2}$  at end, so things only better now.]  
New ball thrown: What is ~~prob.~~ prob. new ball ends up at "height"  $\geq 3$



~~prob.~~ prob. new ball ends up at "height"  $\geq 3$

Both bins it looks at must have  $\geq 2$  balls.

$\therefore$  prob.  $\leq \alpha_2^2 \leq \frac{1}{4}$ .

$\therefore$  "intuitively",  $\alpha_3 =$  fraction of bins with  $\geq 3$  balls  $\leq \frac{1}{4}$   
[each ball thrown has  $\leq \frac{1}{4}$  chance of being in bin of height  $\geq 3$ ]

New ball thrown: what is prob. it ends up at height  $\geq 4$ ? (4)

Both bins it looks at must have  $\geq 3$  balls.

$$\leadsto \text{prob.} \leq \alpha_3^2 \leq \left(\frac{1}{4}\right)^2 \leq \frac{1}{16}$$

$\therefore$  (intuitively)  $\alpha_4 := \text{frac. bins with } \geq 4 \text{ balls} \leq \frac{1}{16}$

$$\alpha_5 \leq \left(\frac{1}{16}\right)^2 = \frac{1}{256} = 2^{-8}$$

$$\alpha_6 \leq \alpha_5^2 \leq 2^{-16}$$

$$\alpha_7 \leq \alpha_6^2 \leq 2^{-32}$$

$$\dots$$
$$\alpha_{k+2} \leq 2^{-2^k}$$

Intuitively if  $\alpha_{k+2} < \frac{1}{n}$ , then probably no bins have  $\geq k+2$  balls

$$\Leftrightarrow 2^{-2^k} < \frac{1}{n} \Leftrightarrow 2^{2^k} > \frac{1}{n} \Leftrightarrow k > \log \log n$$

So "probably" max load  $\leq \log \log n + 2$ .

Ex: 3 choices still "only" gives  $\Theta(\log \log n)$ . [So just stick with 2!]

Bloom Filters: Say:  $\bullet$   $m$  truly enormous (billion, trillion...)

$\bullet$  strings you're storing also large  $\sim L$  bits

(e.g. tweets:  $\sim 280 \times 8 = 2000$  bits)

Can't possibly use less than  $Lm$  bits of space, right? [Right??]

Use  $\approx 8.66m$  bits total!

230x savings [230x fewer servers!!!!] [And hashing with  $n=m$  uses  $\approx Ln$  bits]

What's the catch??

Lookup errors:

Lookup(s) wrongly says "yes" (even tho  $s$  was never stored) w. prob.  $\approx 1.6\%$

$8.66 \approx 1.44k$ ,  $1.6\% = 2^{-k}$ , for  $k = \underline{6}$ . Can choose other  $k$ ...

[trade off: 1.44m extra bits, to halve lookup error prob.]

# How Bloom Filters work

Used in real world:

5

Pick small  $k$ , eg.  $k=6$ .

Google's "BigTable", Hadoop's "HBase", etc.

Set  $N = \frac{k \cdot m}{\ln 2}$  [rounded off to integer], # of bits used.  
 $\frac{1}{\ln 2} \approx 1.44$  [this is the 1.44k bits per item]

Alloc. array  $T[0 \dots N-1]$  of bits, init. all 0's.

Choose  $k$  "independent" hash functions  $h_1, \dots, h_k: \{\text{strings}\} \rightarrow \{0, 1, \dots, N-1\}$

Insert( $s$ ): Set  $T[h_1(s)] = 1, \dots, T[h_k(s)] = 1$ . [maybe some were already 1]

Lookup( $s$ ): Return AND of  $T[h_1(s)], \dots, T[h_k(s)]$ .

Space:  $\approx 1.44k$  bits per item 😊

Time:  $O(1)$  [ $O(k)$ ] operations 😊

Delete( $s$ ): not possible, even slowly 😞

"False positive problem": Lookup( $s$ ) may return True even if  $s$  was never inserted.

[No "false negs" 😊]

Analysis: Prob. [false positive lookup]  $\leq ?$

## Ideas

[again, too fiddly to do rigorously here, so we'll cheat a little]

Q1: After inserting  $m$  items, what frac. of bits in  $T[\cdot]$  do we expect are still?

A: Under SUHA, it's like throwing  $km$  balls into  $N = \frac{k}{\ln 2} m$  bins, (asking about frac. of empty bins)  
 $\lambda = \frac{km}{k/\ln 2 \cdot m} = \ln 2 \approx .69$

$$\Pr[\text{bin } i \text{ empty}] = \left(1 - \frac{1}{N}\right)^{km}$$

Most useful approx. ever:

$$1+x \approx e^x \text{ if } x \text{ is tiny}$$

(6)

$$\hookrightarrow 1+x + \frac{x^2}{2} + \frac{x^3}{6} + \dots$$

If  $|x| = 10^{-6}$   $\frac{10^{-12}}{2} = \text{negligible.}$

∴  $x = -\frac{1}{N}$  is tiny,

$$1 + (-\frac{1}{N}) \approx e^{-\frac{1}{N}}, \quad \therefore \left(1 - \frac{1}{N}\right)^{km} \approx \left(e^{-\frac{1}{N}}\right)^{km} = e^{-\frac{km}{N}}$$

$$(N = \frac{k}{\ln 2} m)$$

$$= e^{-\ln 2} = \frac{1}{e^{\ln 2}} = \frac{1}{2}$$

∴ ~~we~~ Pr [bin is empty]  $\approx \frac{1}{2}$ . [Same for every partic bin.]

∴ We expect after  $m$  inserts,  $T[\cdot]$  is about 50-50 0's and 1's

Now say we do  $\text{Lookup}(s)$ , where  $s$  has never been inserted.  
By SUHA,  $h_1(s), \dots, h_k(s)$  act like  $k$  indep. rand #'s in  $0 \dots N-1$ .

$$\therefore \text{Pr}[\text{AND of } T[h_1(s)], \dots, T[h_k(s)] = \text{True}]$$

$$= \text{Pr}[\text{false pos.}] = \left(\frac{1}{2}\right)^k. \quad \text{☺}$$

Summary: can store  $m$  items of any size

using  $\approx 1.44k$  bits per item,

with  $O(1)$ -time Lookups/Inserts,

and Insert false-positive probs  $\leq 2^{-k}$ .