

Wednesday 9/10/14

## Bloom filter:

huge universe  $U = \{0, 1, \dots, N-1\}$  of size  $N = |U|$ .

0-1 table  $H[0, 1, \dots, n-1]$  of size  $n = |H|$ .

Want to maintain a subset  $S \subset U$ .

$m = |S|$  size of  $S$  after all of the insertions.

Example:  $U =$  possible password strings  
 $S =$  unacceptable passwords

Want fast queries, small space,  
allow false positives with small probability

$k = \#$  of hash functions.

Hash functions  $h_1, h_2, \dots, h_k$  where each  $h_i: U \rightarrow \{0, \dots, n\}$ .

Operations:

- Insert  $x$  into  $S$
- Query is  $x \in S$ ?

Initialize  $H$  to all 0's.

To insert  $x$  into  $S$ :

for all  $i=1 \rightarrow k$ :

[ - Compute  $h_i(x)$   
- Set  $H[h_i(x)] = 1$

To query whether or not  $x \in S$ :

for all  $i=1 \rightarrow k$ :

- Compute  $h_i(x)$

- Check whether  $H[h_i(x)] = 1$ ?

If for all  $i$  it is set = 1, then return (YES)  
else return (NO)

---

For a query: is  $x \in S$ ?

if  $x \in S$ , then we always output YES

if  $x \notin S$ , then we might have a

false positive = incorrectly output YES

What is the probability of a false positive? ③

Recall  $|S|=m$ ,  $|H|=n$ ,

$$\text{let } c = \frac{n}{m} > 1$$

First let's compute the prob. an entry  $i$  is set to 0:

$$\Pr(H[i]=0) = \Pr(\forall y \in S, \forall 1 \leq j \leq k, h_j(y) \neq i)$$

$$= \left(1 - \frac{1}{n}\right)^{km}$$

$$\leq e^{-km/n} = e^{-k/c}$$

because  $e^{-x} = 1 - x + \frac{x^2}{2!} - \frac{x^3}{3!} \pm \dots$

for small  $x$ ,  $e^{-x}$  is a close approx. to  $1-x$ ,  
which is the case for  $\frac{1}{n}$  for large  $n$ .

for  $x \in S$ ,

$$\begin{aligned} & \Pr(\text{false positive}) \\ &= \Pr(\forall_j H[h_j(x)] = 1) \\ &\stackrel{?}{=} (1 - e^{-k/c})^k \end{aligned}$$

This is not correct.

Technically, we should do:

$$\begin{aligned} & \Pr(\forall_j H[h_j(x)] = 1) \\ &= \Pr(H[h_1(x)] = 1) \times \Pr(H[h_2(x)] = 1 \mid H[h_1(x)] = 1) \\ &\quad \times \dots \times \Pr(H[h_k(x)] = 1 \mid H[h_1(x)] = 1, \dots, H[h_{k-1}(x)] = 1) \\ &= 1 - e^{-k/c} \text{ but what about later terms?} \end{aligned}$$

but we do have that:

$$\begin{aligned} \Pr(H[h_2(x)] = 1 \mid H[h_1(x)] = 1) &\leq \Pr(H[h_2(x)] = 1) \\ &\leq 1 - e^{-k/c} \end{aligned}$$

Hence,

$$\Pr(\text{false positive}) \leq (1 - e^{-k/c})^k$$

And we'll use this  $(1 - e^{-k/c})^k$  as our estimate for the false positive rate.

Next question: What's the best choice for  $k$ ? (5)

$k$  big  $\Rightarrow$  putting too many 1's in  $H$

$k$  small  $\Rightarrow$  checking too few bits of  $H$

$$\text{Let } f = (1 - e^{-k/c})^k$$

& let's minimize  $f$  as a function of  $k$ .

$$\text{Let } g = \ln f = k \ln(1 - e^{-k/c})$$

$$\text{so } f = e^g$$

$$\frac{dg}{dk} = \ln(1 - e^{-k/c}) + \frac{k}{1 - e^{-k/c}} \times \frac{1}{c} \times e^{-k/c}$$

$$\text{Set } k = c \ln 2$$

then  $\frac{dg}{dk} = -\ln 2 + \ln 2 = 0$  & can check this is a minimum.

Plugging in  $k = c \ln 2$ , the false positive estimate is

$$f = (1 - e^{-k/c})^k = \left(\frac{1}{2}\right)^{c \ln 2} = \left(\left(\frac{1}{2}\right)^{\ln 2}\right)^c \approx (.6185)^c$$

and  $\Pr(H[i]=0) \approx e^{-k/c} = \frac{1}{2}$

So  $H$  is a random 0-1 string  
(except the bits are not independent  
as discussed earlier)

Examples:

false positive rate

$k=1$ :  $c=10$ : .09516

$c=100$ : .00995

$k=c \ln 2$ :  $c=10$ : .0082

$c=100$ :  $1.3 \times 10^{-21}$

Alternative scheme: Cuckoo hashing

Universe  $U$

static  $S$ : we do a set of insertions and setup  $S$

Then we want to be able to do fast queries on  $S$ .

Goal:  $O(1)$  query time (with no errors)  
& in expectation  $O(1)$  insertion time.

Use 2 hash functions

$$h_1, h_2: U \rightarrow \{0, 1, \dots, n-1\}$$

Store  $\leq 1$  item at each location  $H[i]$ .

To insert  $x$  into  $S$ :

- compute  $h_1(x)$
- if  $H[h_1(x)]$  is empty then  
add  $x$  at  $H[h_1(x)]$

else:

- compute  $h_2(x)$
- if  $H[h_2(x)]$  is empty  
then add  $x$  at  $H[h_2(x)]$   
else let  $y = H[h_2(x)]$   
Set  $H[h_2(x)] = x$   
& move  $y$  to its other  
possible location  
and repeat,

To query: is  $x$  in  $S$ ?

Just check  $H[h_1(x)]$  &  $H[h_2(x)]$ .



For inserting  $x$ ,

this may misplace  $y_1$

moving  $y_1$  may misplace  $y_2$   
etc.

Either this process halts by hitting  
an empty spot or we get a cycle

So if it doesn't halt after  $n$  cycles

then we choose 2 new hash functions  
 $h_1$  &  $h_2$  & we reinsert all of  $S$ .

Pseudocode for insert:

Notation:

$pos$  = position currently trying to insert at

$x \leftrightarrow y$  = swap  $x$  &  $y$

Insert(x):

$pos = h_1(x)$

Do  $\leq n$  times:

[ if  $H[pos] = \emptyset$   
then  $[H[pos] = x$   
     $exit()$   
 $x \leftrightarrow H[pos]$  (So swap  $x$  with item  
    at  $H[pos]$ )  
if  $pos = h_1(x)$  then  $pos = h_2(x)$   
    else  $pos = h_1(x)$

If repeated  $n$  times then Rehash().

Rehash: starts with empty  $H$   
& reinserts every  $x$  into  $S$ .

How often do we rehash?

What's the running time of an insert?

Cuckoo graph:

(11)

Directed graph representing  $H$ .

Vertex for each entry of  $H$ , so  $n$  vertices.

Directed edges show other possible locations for items.

So if  $H[i] = x$  & if  $i = h_1(x)$

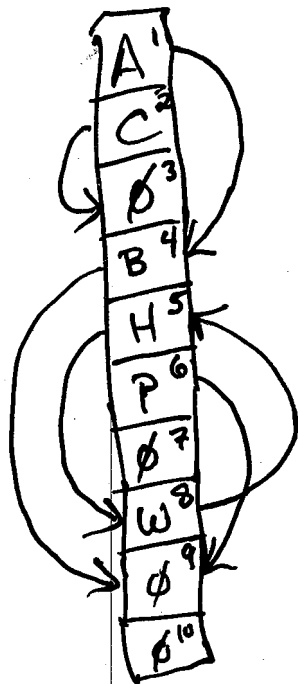
then edge  $i \rightarrow h_2(x)$

if  $H[i] = x$  & if  $i = h_2(x)$

then edge  $i \rightarrow h_1(x)$ .

Example:

Example:



— Suppose we Insert(Z) &  $h_1(Z)=8, h_2(Z)=1$

1. Place Z at pos. 8 & add edge  $8 \rightarrow 1$  (drop edge  $8 \rightarrow 5$ )
2. Move W to pos. 5 & add edge  $5 \rightarrow 8$  (already there)
3. Move H to pos. 8 & add edge  $8 \rightarrow 5$  (drop  $8 \rightarrow 1$ )
4. Move Z to pos. 1 & add edge  $1 \rightarrow 8$  (drop  $1 \rightarrow 4$ )
5. Move A to pos. 4 & add  $4 \rightarrow 1$  (drop  $4 \rightarrow 9$ )
6. Move B to pos. 9 & add  $9 \rightarrow 4$

— Suppose we Insert(D) with  $h_1(D)=5, h_2(D)=8$

Then we get a cycle  $D \rightarrow H \rightarrow W$ :

1. Place D at pos. 5
  2. Move H to pos. 8
  3. Move W to pos. 5
  4. Move D to pos. 8
- etc.

An insertion succeeds if no cycle

If there's a cycle we do a rehash.

First let's ignore rehash's & we'll show the expected insertion time is  $O(1)$ .

Then we'll show that the expected # of rehashes is  $O(1)$ . Since they take  $O(n)$  time per rehash, this is amortized  $O(1)$  time per insert.

$|S|=m$ ,  $|H|=n$ , and we'll choose so that  $n > 6m$ .

Lemma: For  $l \geq 1$ , for positions  $i$  &  $j$ ,  
Probability of a shortest path from  $i \rightsquigarrow j$  of length =  $l \leq \frac{3^{-l}}{n}$

(14)

Hence from the lemma:

Say  $x$  &  $y$  collide if there's a path  $x \rightsquigarrow y$   
or  $y \rightsquigarrow x$ .

In other words, a path from  $\{h_1(x) \text{ or } h_2(x)\}$   
to  $\{h_1(y) \text{ or } h_2(y)\}$   
or from  $\{h_1(y) \text{ or } h_2(y)\}$   
to  $\{h_1(x) \text{ or } h_2(x)\}$

By the lemma,

The Prob. that  $x$  &  $y$  collide is

$$\leq 4 \sum_{l=1}^{\infty} \frac{3^{-l}}{n} = \frac{2}{n} \quad \text{since } \sum_{l=1}^{\infty} 3^{-l} = \frac{1}{2}$$

Hence # of expected collisions with  $x$  is  $O(1)$ .

So when adding  $x$  into  $S$ , there's

$O(1)$  other elements that are moved  
in expectation.

Proof of lemma: Induct on  $l$ .

(15)

Base case:  $l=1$ :

Fix  $i$  &  $j$ . Prob.  $x \in S$  has  $h_1(x)=i, h_2(x)=j$  (or vice versa)

$$= \frac{2}{n^2}$$

Summing over  $x$ ,

Prob. of edge  $i \rightarrow j$  or  $j \rightarrow i$  is

$$\leq m \times \frac{2}{n^2} \leq \frac{1}{3n} \quad \checkmark$$

For  $l > 1$ :

Looking for shortest path of length  $l$  so no length  $l-1$  path.

Take penultimate position  $k$  on the shortest path.

Then (a) there is a path (shortest path) of length  $l-1$

from  $i \rightarrow k$

& (b) there's an edge  $k \rightarrow j$ .

Prob. of (a) by induction is  $\leq \frac{3^{-(l-1)}}{n}$

& (b) by the base case analysis  $= \frac{1}{3n}$

Hence, it's  $\leq \frac{1}{3^{l-1}n} \times \frac{1}{3n} = \frac{1}{3^l n^2}$

Summing over the  $n$  choices of  $k$  we have:

$$\leq \frac{1}{3^l n} \quad \checkmark$$

Rehashing:

To get a rehash we need a cycle.

We'll show that with prob.  $\geq \frac{1}{2}$  no cycles exist.

By lemma,

$$\text{Prob. of a cycle involving position } i \text{ of length } = l \leq \frac{3^{-l}}{n}$$

$$\text{Prob. of some cycle involving } i \text{ is } \leq \frac{1}{n} \sum_{l=1}^{\infty} 3^{-l} = \frac{1}{2n}$$

$$\text{Prob. of some cycle is } \leq n \times \frac{1}{2n} = \frac{1}{2}.$$

Hence, Prob.  $\leq \frac{1}{2}$  of a rehash.

and Prob.  $\leq \left(\frac{1}{2}\right)^k$  of  $k$  rehashes

So expect 1 rehash, and each takes  $O(n)$  time.