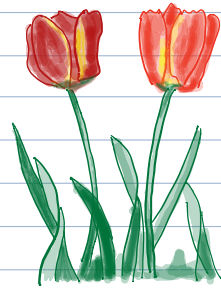


# Secure Distributed Systems

## CompSci 661 / 461

17



These notes:

- Bloom Filters

© 2018-2019 Brian Levine  
All rights reserved.  
Do not distribute or repost.

# Bloom Filters (Harold Bloom, 1970)

A very popular data structure for probabilistically determining if an item  $x$  is an element of a set  $S$ , without transmitting/storing the full set  $S$ .

- the probability of a correct answer is inversely proportional to the size of the compact representation of  $S$ .

## Some preliminary definitions:

The Bloom Filter has a chance of answering incorrectly. To evaluate/characterize its performance, we need to distinguish the ground truth from the result of asking the Bloom Filter.

$$\text{ground truth} \begin{cases} \text{Positive} - x \in S & (x \text{ is in } S) \\ \text{Negative} - x \notin S & (x \text{ is not in } S) \end{cases}$$

A guess is True if the guess is correct;  
and False if the guess is incorrect.

Hence, there are four possibilities

	Positive	Negative
True	TP	TN
False	FP	FN

Eg. FP is an incorrect guess that  $x \in S$ , when in fact  $x \notin S$ .

The False Positive rate (FPR) is

$$\frac{FP}{FP + TN} \quad \leftarrow \begin{array}{l} \text{number of False Positives} \\ \text{number of Negatives (all of them)} \end{array}$$

The False Negative Rate (FNR) is

$$\frac{FN}{TP + FN} \quad \leftarrow \begin{array}{l} \text{number of False Negatives} \\ \text{number of Positives (all of them)} \end{array}$$

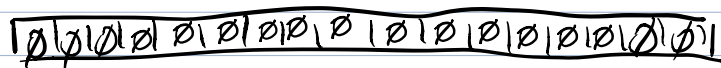
Bloom Filters have an

- Adjustable False Positive Rate (FPR)
- Zero False Negative Rate (FNR)

Let's see what's involved

## Here's how Bloom Filters work.

- ① Create an array of  $m$  bits, all cleared (zeros)



- ② For each element  $i$  of  $S$ , send it thru a hash function  $K$  times with prefixes  $x_1, x_2, \dots, x_K$

$$h(x_1, i) \bmod m = h_1$$

$$h(x_2, i) \bmod m = h_2$$

$\vdots$

- we have  $K$  different hashes of the same input  $i$
- $h_1$  is between 0 and  $(m-1)$ , of course.
- Set indices  $h_1, h_2, \dots, h_K$  in the bit array.
- IF they are already set, keep them set.

- ③ To test if  $j \in S$ : send  $j$  thru the  $K$  hashes.

$$h(x_1, j) \bmod m = h_1 \text{ etc}$$

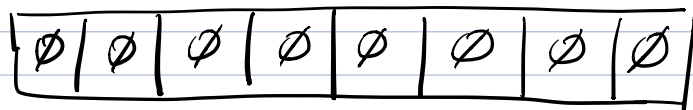
Check each index in the array ( $h_1, h_2, \dots, h_K$ )

- IF any are 0, there is zero chance  $j$  is in  $S$ .
- IF all are 1, we report that  $j$  is in  $S$ , but there is a non-zero chance we are wrong (False Positive)

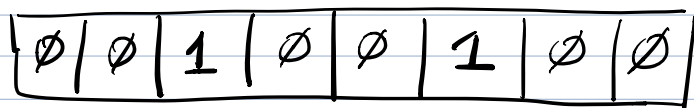
In sum  $n$  items,  $m$  bits in array,  $K$  hash functions

## Example.

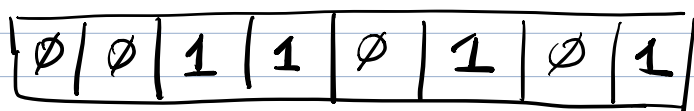
Set  $n = 16$  bits in array  
 $k = 2$   
 $b = 4$



insert 1 element "a". And say  $h_1 = \text{hash}(x_1, \text{element}) \otimes 16 = 2$   
 $h_2 = \text{hash}(x_2, \text{element}) \otimes 16 = 5$



We might insert a few more elements ...



Now check if "a" is in our filter.

$\left. \begin{array}{l} h_1 = 2 \\ h_2 = 5 \end{array} \right\}$  both are 1, therefore is (supposedly)

If we check for "b" and find that

$\left. \begin{array}{l} h_1 = 2 \\ h_2 = 1 \end{array} \right\}$  Since at least one index is  $\emptyset$ , we know that "b" is not in the original set.

how do we get a false positive? Say we never inserted "c" into the filter. But it just so happens that

$\left. \begin{array}{l} h_1 = 2 \\ h_2 = 3 \end{array} \right\}$  Since both are unfortunately set, we have a false positive.

The false positive rate depends on the number of elements — not the length of each element.

## What is the FPR of a Bloom Filter?

(from Wikipedia)

Let's derive the false positive rate.

- For each hash function, we set 1 index in the array

The probability that one of the  $m$  bits is set is  $(1/m)$

- Probability of not being set is  $(1 - 1/m)$

- Prob of not being set by  $K$  hash functions:  $(1 - 1/m)^K$

- Let  $n = |S|$ . Since we will insert  $n$  elements, the probability that an index of the  $m$  bit array is not set is  $(1 - 1/m)^{Kn}$

- Probability that an index is set is:  $1 - (1 - 1/m)^{Kn}$

The Probability that  $K$  indices are set for an element that is not in  $S$  is:  $(1 - (1 - 1/m)^{Kn})^K$

which is approximately:

$$P \approx \left(1 - e^{-\frac{Kn}{m}}\right)^K$$

## How many hash functions should we use?

It has been shown that  $K$  is optimal for a desired FPR of  $p$  when  $K = -\log_2(p)$

$$\text{E.g., for } p = 0.01 \Rightarrow K = \lceil -\log_2(0.01) \rceil = \lceil 6.6 \rceil = 7$$

How many hash functions should we use?

[long version]

We want the value of  $K$  that minimizes the FPR

From  
Mitzenmacher

So we take the derivative of the FPR w.r.t.  $K$

$$\begin{aligned} p &= (1 - e^{-\frac{kn}{m}})^k \\ &= (1 - \exp(-\frac{kn}{m}))^k \\ &= \exp(\ln(1 - \exp(-\frac{kn}{m})))^k \\ &= \exp(k \ln(1 - \exp(-\frac{kn}{m}))) \end{aligned}$$

let  $e^x = \exp(x)$  for clarity

since  $a = \exp(\ln(a))$

Minimizing  $x$  also minimizes  $\exp(x)$

$$\text{let } g = k \ln(1 - e^{-kn/m})$$

$$\frac{\partial g}{\partial k} = \ln(1 - e^{-kn/m}) + \frac{kn}{m} \cdot \frac{e^{-kn/m}}{1 - e^{-kn/m}}$$

This derivative is equal to 0 when  $k = \ln(2) \frac{m}{n}$  (not shown here)

Let's substitute our optimal value of  $k$  into our equation for the FPR, and solve for  $m$

$$P = \left(1 - e^{-\frac{kn}{m}}\right)^k$$

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

$$= \left(1 - e^{-\ln(2) \frac{m}{n} \frac{n}{m}}\right)^{\ln(2) \frac{m}{n}}$$

$$= \left(1 - e^{-\ln 2}\right)^{\ln(2) \frac{m}{n}}$$

$$= \left(1 - \frac{1}{2}\right)^{\ln(2) \frac{m}{n}}$$

$$P = \left(\frac{1}{2}\right)^{\ln(2) \frac{m}{n}}$$

$$\ln(P) = \ln\left(\left(\frac{1}{2}\right)^{\ln(2) \frac{m}{n}}\right)$$

$$= \ln(2) \frac{m}{n} \ln\left(\frac{1}{2}\right)$$

$$= \frac{m}{n} \ln(2) (-\ln(2))$$

$$\ln(P) = -\frac{m}{n} \ln^2(2)$$

$$\ln^2(2) = \ln(2) \ln(2)$$

$$\frac{n \ln P}{\ln^2(2)} = -m$$

$$\ln^2(2)$$

$$m = \frac{-n \ln P}{\ln^2(2)}$$

← number of bits in array given desired FPR of  $P$  and  $n$ .

And given  $n$  and  $m$ , the optimal number of hash functions is  $k = \frac{m}{n} \ln 2$ , (not an integer!)

Note that  $\frac{m}{n} = \frac{-\ln P}{\ln^2(2)}$

or

and so  $k = \frac{-\ln P}{\ln^2(2)} \cdot \ln 2 = \frac{-\ln P}{\ln 2} \Rightarrow k = -\log_2 P$