



Design and Analysis  
of Algorithms I

# Data Structures

---

## Bloom Filters

# Bloom Filters: Supported Operations

Raison D'être: fast Inserts and Lookups.

Comparison to Hash Tables:

Pros: more space efficient.

Cons:

- 1) can't store an associated object
- 2) No deletions
- 3) Small false positive probability

(i.e., might say x has been inserted even though it hasn't been)

# Bloom Filters: Applications

Original: early spellcheckers.

Canonical: list of forbidden passwords

Modern: network routers.

- Limited memory, need to be super-fast

# Bloom Filter: Under the Hood

Ingredients: 1) array of  $n$  bits ( *So  $\frac{n}{|S|}$  = # of bits per object in data set  $S$*  )

2)  $k$  hash functions  $h_1, \dots, h_k$  ( $k$  = small constant)

Insert( $x$ ): for  $i = 1, 2, \dots, k$  (whether or not bit already set to 1)  
set  $A[h_i(x)] = 1$

Lookup( $x$ ): return TRUE  $\Leftrightarrow A[h_i(x)] = 1$  for every  $i = 1, 2, \dots, k$ .

Note: no false negatives. (if  $x$  was inserted, Lookup ( $x$ ) guaranteed to succeed)

But: false positive if all  $k$   $h_i(x)$ 's already set to 1 by other insertions.

# Heuristic Analysis

Intuition: should be a trade-off between space and error (false positive) probability.

Assume: [not justified] all  $h_i(x)$ 's uniformly random and independent (across different  $i$ 's and  $x$ 's).

Setup:  $n$  bits, insert data set  $S$  into bloom filter.

Note: for each bit of  $A$ , the probability it's been set to 1 is (under above assumption):

Under the heuristic assumption, what is the probability that a given bit of the bloom filter (the first bit, say) has been set to 1 after the data set  $S$  has been inserted?

☐  $(1 - 1/n)^{k|S|}$  prob 1<sup>st</sup> bit = 0

☐  $1 - (1 - 1/n)^{k|S|}$  prob 1<sup>st</sup> bit = 1

☐  $(1/n)^{|S|}$

☐  $(1 - 1/n)^{|S|}$

# Heuristic Analysis

Intuition: should be a trade-off between space and error (false positive) probability.

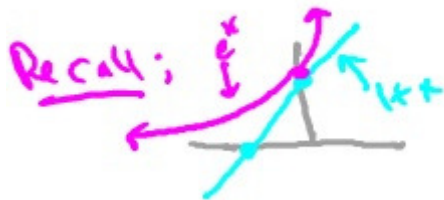
Assume: [not justified] all  $h_i(x)$ 's uniformly random and independent (across different  $i$ 's and  $x$ 's).

Setup:  $n$  bits, insert data set  $S$  into bloom filter.

Note: for each bit of  $A$ , the probability it's been set to 1 is (under above assumption):

$$1 - \left(1 - \frac{1}{n}\right)^{k|S|} \leq 1 - e^{-\frac{k|S|}{n}} = 1 - e^{-\frac{k}{b}}$$

$b = \#$  of  
bits per  
object  
( $n/|S|$ )



# Heuristic Analysis (con'd)

Story so far: probability a given bit is 1 is  $\leq 1 - e^{-\frac{k}{b}}$

So: under assumption, for x not in S, false positive probability is  $\leq [1 - e^{-\frac{k}{b}}]^k$   
where b = # of bits per object.

How to set k?: for fixed b,  $\epsilon$  is minimized by setting

Plugging back in:

$$\epsilon \approx \left(\frac{1}{2}\right)^{(\ln 2)b} \quad \text{or} \quad b \approx 1.44 \log_2 \frac{1}{\epsilon}$$

(exponentially  
small in b)

$$k \approx (\ln 2) \cdot b \approx 0.693$$

error rank  $\epsilon$

Ex: with b = 8, choose k = 5 or 6, error probability only approximately 2%.