

BF = Bloom filters

set S of n keys : membership $x \in S$ 1-side error with pr. f

$$f \sim (1-p)^k \text{ failure probability}$$

$$\left. \begin{aligned} k &= \frac{m}{n} \ln 2 \quad \# \text{ hash function} \\ p &= \frac{1}{2} \quad \text{best } p = \frac{1}{2} \end{aligned} \right\} \Rightarrow f \sim \frac{1}{2^{\frac{m}{n} \ln 2}} \sim 0.618$$

space : m bits + $O(k)$ words

chosen to get a certain f

More info on the space:

$$\log_2 f \sim \frac{m}{n} \ln 2 \cdot \lg \frac{1}{2} \quad \boxed{-1}$$

$$\Rightarrow \boxed{m} \sim \frac{n \lg f}{\lg \frac{1}{2} \ln 2} = \frac{n \lg \frac{1}{f}}{\ln 2} \sim \boxed{1.44 \cdot n \lg_2 \frac{1}{f}}$$

$\# \text{ bits}$

Approximate dictionary

- set S of n keys
- $x \in S$ with error probability f

BF is one possible answer
Is it the best possible?

① LOWER BOUNDS

- TRIVIAL: $\Omega(1)$ time
- $\geq n \lg_2 \frac{1}{f}$ bits

$O(k)$ time

space $\sim 1.44 \lg_2 \frac{1}{f}$
bits per key

② UPPER BOUNDS

- $O(1)$ time with 1 hash function
- $\sim n \lg_2 \frac{1}{f}$ bits + lower order terms

Information theory lower bound

$S \subseteq U$ and let $|U| = m$ ($\neq m$ in BF!)

$$|S| = n$$

Q. How many sets $S \subseteq U$ of size n ? $\binom{m}{n}$ choices

I.T. says that using less than $\lg_2 \binom{m}{n}$ in the worst case cannot give a correct algorithm: indeed, using less bits forces two S' and S'' to get the same binary representation.

so membership is forced

exact

A useful formula:

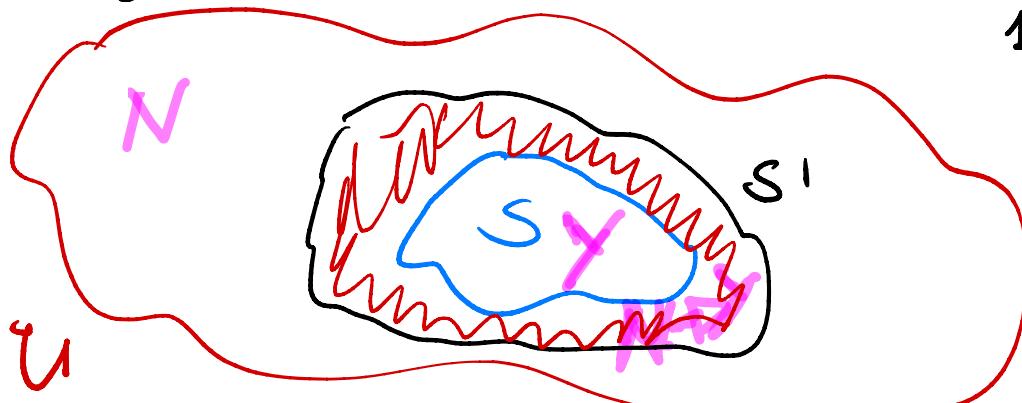
$$\lg_2 \binom{m}{n} \sim n \boxed{\lg_2 \frac{m}{n}}$$

Approximate dictionary for S with pr. f , D' $\{$ exact dictionary for S'

Exact dictionary for some $S' \supseteq S$ and $\frac{|S' \setminus S|}{|U|} = f$

$S' = \{x \in U : \text{approximate dict says YES}\}$

$S' = \{x \in U : D'(x) = \text{True}\}$



1-side error for $S \Rightarrow$

$S \subseteq S' \Rightarrow$ all keys in S are accepted plus the extra (wrong) keys in $S' \setminus S$

Q: Let D' be the exact dictionary for S'

D " " " S

(note that D' is also the approximate dictionary for S)

Given D' , can we get D ?

$b' = \# \text{bits required by } D'$

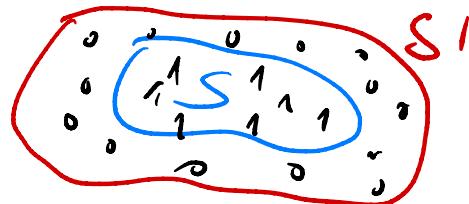
$D = D' + \text{extra bits to mark correct answers from } D'$

we have $|S'|$ YES answers, and we mark

the i -th YES answer from D'

$\begin{cases} 1 & x \in S \\ 0 & x \notin S \end{cases}$

(have
 $x \in S \setminus S'$)



We get an exact dictionary $D \Rightarrow$ it must require $\log_2(m)$
bits by I.T.

$$|D| \geq \log_2(m) \quad \text{I.T.}$$

$|D'| + \text{extra bits} \geq \log_2(m)$ as $D' + \text{extra bits}$ are a
dictionary for S'

$$|D'| + \log_2(|S'|) = b' + \log_2(n + fm)$$

$$(|S'| = |S| + |S'/S| = |S| + |U| \cdot f = n + fm)$$

$$\frac{|S'/S|}{|U|} = f$$

\uparrow
 $|S'/S| = fm$

Putting these parts together :

$$b' + g\left(\frac{n+fm}{n}\right) > g\left(\frac{m}{n}\right)$$

$$\begin{aligned} b' &\geq g\left(\frac{m}{n}\right) - g\left(\frac{n+fm}{n}\right) \sim n g \frac{m}{n} - n g \underbrace{\frac{n+fm}{n}}_{\sim g \frac{fm}{n}} \sim n g \frac{m}{n} - n g \frac{fm}{n} \\ &= n \left(g \frac{m}{n} - g \frac{fm}{n} \right) \\ &= n \left(g \frac{m}{n} + g \frac{n}{fm} \right) \\ &= n \left(g \frac{m}{n} \cdot \cancel{\frac{n}{fm}}^1 \right) = n g_2 \frac{1}{f} \end{aligned}$$

$$\sim g \frac{fm}{n} \Rightarrow \frac{n+fm}{n} = 1 + \frac{fm}{n}$$

D' requires at least $n g_2 \frac{1}{f}$

UPPER BOUND

Bit vector m bits n 1s

[RRR] $\binom{m}{n} + \text{lower order terms}$

constant-time lookup
succinct data structure

$B(i)$ in $O(1)$ time

①

Universal Hashing $h \in \mathcal{H}$

$$\frac{1}{k} = \frac{m}{n} \Rightarrow m = \frac{n}{k}$$

②

$\forall x \in S : B[h(x)] = 1$ (BF with 1 hash function
RRR⁺)

- Space is optimal because of RRR
- time is $O(1)$ as we use 1 hash function
- error probability : union bound

$$\Pr[B[i] = 1] = \sum_{x \in S} \Pr[h(x) = i] = \underbrace{\frac{n}{m}}_{\frac{1}{m}} = f$$

$$m = \text{universe} = |U| = \frac{n}{f}$$

you store the hash values in B^S