

Bloom filters

Recall:

$$\begin{matrix} \rightarrow |X| = n \\ X \subseteq U \end{matrix}$$

- represents set X in low memory
- $\text{Insert}(x)$
- $\text{Search}(x) \begin{cases} x \in X \Rightarrow \text{always True} \\ x \notin X \Rightarrow \text{sometimes also True} \\ \text{(false positive)} \end{cases}$

• implementation:

- bit array $B[0 \dots m-1]$
- hash function $h: U \rightarrow [m]$
- $\text{Insert}(x) \rightarrow B[h(x)] \leftarrow 1$
- $\text{Search}(x) \rightarrow \text{return } B[h(x)]$

single
band
filter

• $\text{Pr}[\text{false positive}] \leq \frac{n}{m}$

• multiband filter = k independent filters

$(B_i, h_i) = \text{band}$

$\rightarrow \text{Pr}[\text{FP}] \leq \left(\frac{n}{m}\right)^k = \frac{1}{2^k} \text{ for } m=2n$

also possible to use
single table of $k \cdot m$ bits
and k functions $h_i: U \rightarrow [k \cdot m]$

Optimizing parameters

• assume n given

• $\epsilon > 0$ given error probability

Thm: Set $m=2n$, $k = \lceil \log(1/\epsilon) \rceil$.

Then we get $\text{Pr}[\text{FP}] \leq \epsilon$.

\hookrightarrow need $2n \cdot \lceil \log(1/\epsilon) \rceil$ bits

\hookrightarrow for $\epsilon = 0.01 \approx 14 \text{ Mb}$

vs $\approx 100 \text{ Mb}$ for single band

• can we do better memory-wise?

\hookrightarrow can we choose better m and k ?

- given M bits
 → find m, k s.t. $m \cdot k \leq M$
 with ϵ smallest
- assume perfectly random hash functions

• single band in the filter:

$$\rightarrow \Pr[B_+^i = 0] = \underbrace{\left(1 - \frac{1}{m}\right)^n}_{\Pr[h_\epsilon(x_j) \neq i]}$$

m is parameter!

$$\rightarrow \left(1 - \frac{1}{m}\right)^n \approx e^{-\frac{n}{m}} = p$$

↳ \leq and \approx for $n \rightarrow \infty$

better bound for $\Pr[FP]$ than $\frac{n}{m}$ from union bound (but needs perfectly random functions)

⊙ value of p determines other parameters

$$\rightarrow m \approx -n / \ln p$$

↳ base e log!

$$\rightarrow k = \lfloor M/m \rfloor \approx -M/n \cdot \ln p$$

↳ bands must fit into M bits

$$\rightarrow \epsilon = \Pr[FP] = \Pr[\exists i B_+^i = 1]$$

$$= (1-p)^k = e^{k \cdot \ln(1-p)} \approx \exp\left(\underbrace{-M/n}_{\text{fixed}} \cdot \underbrace{\ln(1-p) \cdot \ln(p)}_{\text{maximize w.r.t. } p \text{ to get min } \epsilon}\right)$$

⊙ $\min_k \epsilon$ is for $p = 1/2$

↳ maximize $\ln(1-p) \cdot \ln(p)$ s.t. $p \in (0, 1)$ using basic calculus

$$\Rightarrow \epsilon = 2^{-k} = 2^{-M/n \cdot \ln 2}$$

⊙ For fixed ϵ and minimizing M we get:

$$\rightarrow \text{fix } m = n / \ln 2$$

$$\rightarrow k = \lceil \log(1/\epsilon) \rceil \quad (\text{to get } \Pr[FP] \leq \epsilon)$$

$$\rightarrow M = k \cdot m \approx n \cdot \log(1/\epsilon) \cdot (1/\ln 2) = n \cdot \log(1/\epsilon) \cdot 1.44$$

↳ vs $n \cdot \log(1/\epsilon) \cdot 2$ for $m = 2n$

Note: It can be shown that no data structure can do better than $n \cdot \log(1/\epsilon)$ for this problem

Counting filters

? How to implement $Delete(x)$?

→ setting $B[h(x)] \leftarrow 0$ could cause false negatives if $\exists y \in X$ s.t. $h(y) = h(x)$

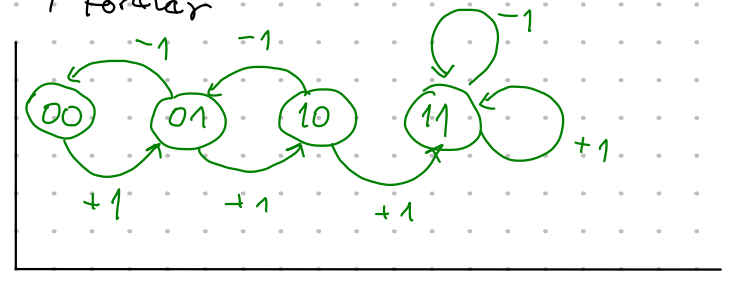
Solution: use counters

- $C[0 \dots m-1]$ array of b -bit counter
 - $Insert(x) \rightarrow C[h(x)]++$
 - $Delete(x) \rightarrow C[h(x)]--$
 - $Search(x) \rightarrow$ return $C[h(x)] > 0$
- still allows false positives
- $\exists y \in X$ s.t. $h(y) = t \Rightarrow C[t] > 0$
→ no false negatives

Problem: what if counter overflows?

Solution: counter gets "stuck" at $2^b - 1$ forever

① "simple" filter is counting filter with $b=1$



? What is the probability that $C[i]$ gets stuck?

- assume single band and fully random hash function, after n inserts.
- $C[i] = \#$ increments of $C[i]$ (e.g. what would be in $C[i]$ if we had unlimited precision)

$$Pr[C[i] = t] = \binom{n}{t} \cdot \left(\frac{1}{m}\right)^t \cdot \left(1 - \frac{1}{m}\right)^{n-t}$$

- for each of x_j we toss a coin
- with probability $1/m$ x_j falls into bucket i
- binomial distribution with $p = 1/m$

$$\Rightarrow Pr[C[i] \geq t] \leq \binom{n}{t} \cdot \left(\frac{1}{m}\right)^t$$

→ we don't care whether rest of the elements end up in i or not

$$\rightarrow t = 2^b - 1$$

$$P_r[C[i] \text{ is stuck}] \leq P_r[C'[i] \geq t] \leq \left(\frac{n \cdot e}{t}\right)^t \cdot \left(\frac{1}{m}\right)^t$$

\hookrightarrow for any t : $\binom{n}{t} \leq \left(\frac{ne}{t}\right)^t$

\rightarrow for $m = n/\ln 2$:

$$P_r[C[i] \text{ is stuck}] \leq \left(\frac{e \cdot \ln 2}{t}\right)^t \rightarrow \text{for } b=4 \rightsquigarrow \leq 10^{-13}$$

$$P_r[\exists \text{st. } C[i] \text{ is stuck}] \leq m \cdot \left(\frac{e \cdot \ln 2}{t}\right)^t \rightarrow \text{for } m=10^9 \rightsquigarrow \leq 10^{-4}$$

Note: We can periodically rebuild filter to get rid of stuck filters
 \hookrightarrow but we need to know "true" content of X