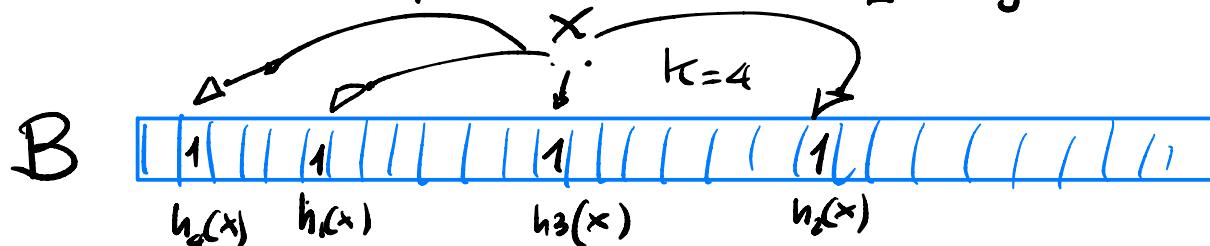


BLOOM FILTERS: it does not store the keys (MonteCarlo algs.)

- set $S \subseteq U$ of keys (but we do not want to store them in the BF!)
- universal hash family \mathcal{H}
- k m.s. $h_1, h_2, \dots, h_k \in \mathcal{H}$ randomly and uniformly chosen
 $h_i : U \rightarrow [m]$

BUILDING

- $B = \text{bitvector} = \text{compact array of } m \text{ bits}$ (initialized to all 0s) (like C++ bitset)
- Ideas: $x \in S \xrightarrow{\text{BF}}$ for $i=1, 2, \dots, k$: $B[h_i(x)] = 1$



What we store:

e.g. $m, k, p, n = |S|$, etc.

► B m bits + $O(1)$ words

► h_1, h_2, \dots, h_k $2k$ words

Space cost is dominated by m . Q. How to choose m, k ?

Implicit condition:

$$x \in S \Leftrightarrow B[h_1(x)] = 1 \wedge B[h_2(x)] = 1 \wedge \dots \wedge B[h_k(x)] = 1$$



Error due to the implicit condition:

1-side error:

if $x \in S$ really \Rightarrow $\textcircled{*}$ is true : however, $\textcircled{*} = \text{true}$ \neq

if $x \notin S$ \Rightarrow $\textcircled{*}$ is false : it is possible that $x \notin S$ FAILURE PROBABILITY

Query: x vs test \otimes in $O(k)$ time

Insertion: x vs set \otimes in $O(k)$ time

Deletion: it cannot be executed in this version
(use counters)

- Q:
- ① error probability f
 - ② choose "best" parameters

- Universal hashing: $\Pr[h_i(x) = y] = \frac{1}{m}$
- $f = \Pr(\text{ERROR}) = \Pr(x \notin S \text{ but } h_i(x) = \text{true})$

Small steps: suppose we built B , take any position $q \in [m]$

$$\underline{Q_1}: \Pr[B[q] = 1] = 1 - \Pr[B[q] = 0]$$

$$\bar{Q}_1: \text{fix } h_i : \Pr[B[q] = 0 \mid h_i \text{ is chosen}] = 1 - \frac{1}{m}$$

as $\Pr[h_i(x) = q] = \frac{1}{m}$

$$\bar{Q}_1: \Pr[B[q] = 0] = \left(1 - \frac{1}{m}\right)^k \quad \text{as all } h_i \text{'s must give values } \neq q$$

Repeating the above argument for all keys in S

$$\Pr[B[q] = 0] = \left(\left(1 - \frac{1}{m}\right)^k\right)^n \quad \text{: this is what we have after building } B$$

$$P' = \left(1 - \frac{1}{m}\right)^{nk} \quad \text{and} \quad B[q] = 0$$

$1 - P'$ $B[q] = 1$



$$f = \Pr[\otimes] = (1 - p')^k \quad \text{where } p' = \left(1 - \frac{1}{m}\right)^{nk}$$

FAILURE PROBABILITY that we want to minimize

Taylor expansion $(1 + y) \approx e^y$ for small y

$$P' = \underbrace{\left(1 - \frac{1}{m}\right)^{nk}}_y \approx \left(e^{-\frac{1}{m}}\right)^{nk} = e^{-\frac{nk}{m}} = P$$

We minimize $(1 - P)^k$

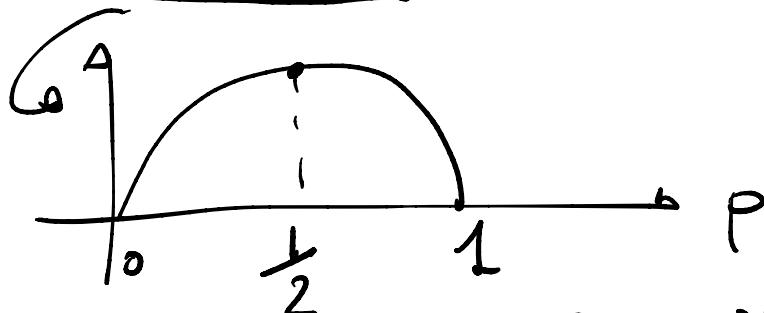
minimize $f \sim (1-p)^k$ where $p = e^{-\frac{nk}{m}}$

take logs. $\ln(1-p)^k = k \ln(1-p)$

$$\ln p = -\frac{nk}{m} \cancel{\ln e} \Rightarrow k = -\frac{m}{n} \ln p$$

replace k

$$k \ln(1-p) = -\frac{m}{n} \cdot \ln p \cdot \ln(1-p) \leftarrow \text{minimize over } p$$



Best choice is $p = \frac{1}{2}$

$$f \sim (1-p)^k = \frac{1}{2^k} = \frac{1}{2^{-\frac{m}{n} \ln \frac{1}{2}}} = \frac{1}{2^{\frac{m}{n} \ln 2}} = \underbrace{\left(2^{-\ln 2}\right)^{m/n}}_{\approx 0.618} < 0.618^{\frac{m}{n}}$$

$$k = -\frac{m}{n} \ln p = \frac{m}{n} \ln 2$$

You have to choose m accordingly