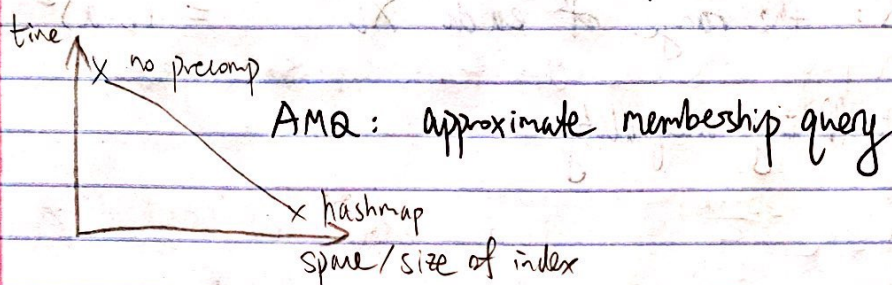2/11/21     Hash-Based Data Structures

Membership query

① no precomputation is allowed     $O(N)$ time

② precomputation is allowed : hashmap   $O(1)$ time, $O(N)$ space



AMQ : approximate membership query

AMQ     Data Structure : Bloom Filter

Data item $S$

Hash functions $h_i$ : $\{h_1, \cdots, h_K\}$

$h$ : set elm $\rightarrow \{0, \cdots, m-1\}$

item $s \rightarrow$ signature $[h_1(s), h_2(s), \cdots, h_K(s)]$     func

binary vector of size $m$ indicating that at least 1 hash ~~func~~

has that value

$sig(s) \rightarrow$ $m$-dim binary vector   has at most $k$ ones

A set of items $S = \underset{N}{=} \{s_1, \cdots, s_N\}$

$Sig(S) = \underset{i=1}{\text{BitWise OR}} \{sig(s_i)\}$

Ex.     $i \rightarrow$ new item ; $Sig(S)$ ; $sig(i)$

If $sig(\overset{i}{\mathbf{\imath}})$ is 1 somewhere where $Sig(S)$ is 0

$\Rightarrow$     $S$ does not contain $i$

$Sig(S)$ can reject values not in the set

If $sig(i)$ is 0 somewhere where $sig(S)$ is 1

$\Rightarrow$     $S$ might contain $i$

The 1 could've been set by another item

Scanned with CamScanner

If $Sig(i)$ has a 1 and $Sig(s)$ has a 1 in the same spot
$\Rightarrow$ $S$ might not contain $i$
$\vdots$ could've been set by some other elements

BitwiseOR sigs : no false negatives, possibly false positives
usage : partition

Analysis $\epsilon = (1 - e^{-\frac{kN}{m}})^k$ ← # of hash funcs
False positive rate & $k$ & $m$ ← size of bloom filter

vs $O(N)$ hashmap
$O(m)$ space

① Each hash func : $\frac{1}{m}$ prob. to set a bit to 1
② Over $k$ applications to a single data point, prob. of a given bit not 1 ?
$P = (1 - \frac{1}{m})^k$
③ $e^{-1} = \lim_{m \to \infty} (1 - \frac{1}{m})^m$ $\Rightarrow$ $p = (1 - \frac{1}{m})^k \approx e^{-\frac{k}{m}}$
(ex. in database)
④ $N$ datapoints : $p^N = e^{-\frac{kN}{m}}$
over $k$ applications to $N$ data points, prob. that a given bit not 1

⑤ A False Positive happens when all bits of a new item are 1
$(1 - P^N)^k = \boxed{(1 - e^{-\frac{kN}{m}})^k} = \epsilon$ : false positive rate
space of data struct

Q  How do we set $m$ and $k$ ?        choose $k$ to minimize FPR
$\epsilon = \gamma^k$ , $\gamma = (1 - e^{-\frac{kN}{m}})$ , want $\gamma$ to have high entropy, i.e. $\gamma = \frac{1}{2}$
$\Rightarrow$ $k = \frac{m}{N} \ln 2$

Q  How to use a signature to estimate the number of input items ?
$\bar{E}(\hat{r}) = $ frac of bits that are 1
① $N = \dfrac{-m \log(1 - \bar{E}(\hat{r}))}{k}$