

Problem Given set S of n numbers in $\{0, 1, \dots, U-1\}$, build data structure st. we can quickly answer membership queries: given y , is $y \in S$?

("Dictionary")

e.g. sorted array

\Rightarrow $O(\log n)$ query time
 $O(n \log n)$ preprocessing time
 $O(n)$ space

balanced search tree

\Rightarrow $O(\log n)$ insert/delete time
 \leftarrow "dynamic"

("predecessor/successor search")

better?

Method 0.

bit vector



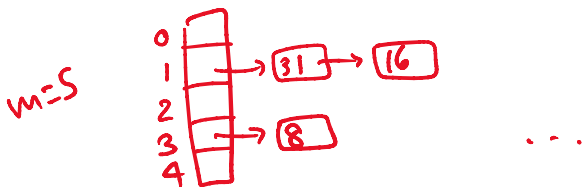
$O(1)$ time
 $O(U)$ space

Method 1. hash table with chaining

Pick a hash fn $h: \{0, \dots, U-1\} \rightarrow \{0, \dots, m-1\}$ for some $m \ll U$.

Store array $A[0, \dots, m-1]$ where $A[i] =$ list of all $x \in S$ with $h(x) = i$
 \leftarrow "bucket"

e.g. $h(x) = x \bmod m$



query(y):
search the list $A[h(y)]$

→ if input is rand, unif. distributed,
each bucket will have $\sim \frac{n}{m}$ elems
"on average"

Set $m \approx n \Rightarrow O(n)$ space
 $O(1)$ average query time
 $O(n)$ preproc. time

but can't assume input is random!

Method 2 hashing with random hash fn
chosen a universal family
[Carter, Wegman '77]

Ex Assume U is prime.

Pick rand. $a \in \{1, \dots, U-1\} \leftarrow$
rand $b \in \{0, \dots, U-1\} \leftarrow$

Define $h_{a,b} : \{0, \dots, U-1\} \rightarrow \{0, \dots, m-1\}$:

$$h_{a,b}(x) = \left((ax + b) \bmod U \right) \bmod m$$

computable in $O(1)$ time

Property For any fixed $x, y \in \{0, \dots, U-1\}$ ($x \neq y$),

(called universality) → $\Pr_{a,b} \left[h_{a,b}(x) = h_{a,b}(y) \right] \leq \underline{\underline{O\left(\frac{1}{m}\right)}}$.

↑
we say x, y collide

More strongly, for any fixed $x, y \in \{0, \dots, U-1\}$, ($x \neq y$)
 $i, j \in \{0, \dots, m-1\}$,

(called strong 2-universality) → $\Pr_{a,b} \left[h_{a,b}(x) = i \wedge h_{a,b}(y) = j \right] \leq O\left(\frac{1}{m^2}\right)$.

Pf: Fix $i, j \in \{0, \dots, U-1\}$.

Pf: Fix $i, j' \in \{0, \dots, U-1\}$.

$$\Pr_{a,b} \left[\begin{array}{l} (ax+b) \bmod U = i' \\ \wedge (ay+b) \bmod U = j' \end{array} \right]$$

$$= \Pr_{a,b} \left[\begin{array}{l} ax+b \equiv i' \pmod{U} \\ \wedge ay+b \equiv j' \pmod{U} \end{array} \right]$$

$\left. \begin{array}{l} \text{2x2 linear system of eqns} \\ \text{soln is uniq.} \end{array} \right\}$

$$= \frac{1}{U(U-1)}$$

$$\left\{ \begin{array}{l} a = (i' - j') (x - y)^{-1} \\ b = i' - ax \end{array} \right.$$

$$\Rightarrow \Pr_{a,b} [h_{a,b}(x) = i \wedge h_{a,b}(y) = j]$$

$$= \sum_{\substack{i', j' \in \{0, \dots, U-1\} \\ i' \bmod m = i \\ j' \bmod m = j}} \frac{1}{U(U-1)}$$

$$\leq \lceil U/m \rceil \lceil U/m \rceil \cdot \frac{1}{U(U-1)}$$

$$= O\left(\frac{1}{m^2}\right). \quad \square$$

Query time analysis:

for a fixed query y ,

$$E(\text{query-time}) \approx E(\# \text{ collisions with } y)$$

$$= E\left[\left| \{x \in S - \{y\} : h_{a,b}(x) = h_{a,b}(y)\} \right| \right]$$

$$= E\left[\sum_{x \in S - \{y\}} \mathbb{1}_{[h_{a,b}(x) = h_{a,b}(y)]} \right]$$

(indicator fn: $1_E = \begin{cases} 1 & \text{if } E \text{ true} \\ 0 & \text{else} \end{cases}$)

$$(E[1_E] = 1 \cdot \Pr(E) + 0 \cdot \Pr(E^c))$$

$$= \sum_{x \in S - \{y\}} E \left[\underline{1_{[h_{a,b}(x) = h_{a,b}(y)]}} \right]$$

by linearity of expectation

$$= \sum_{x \in S - \{y\}} \Pr[h_{a,b}(x) = h_{a,b}(y)]$$

$$\leq \sum_{x \in S - \{y\}} O\left(\frac{1}{m}\right) \quad \text{by universality}$$

$$= O\left(\frac{n}{m}\right)$$

Set $m \approx n \Rightarrow$

- $O(n)$ space
- $O(1)$ expected query time
- $O(n)$ preproc time
- $O(1)$ insert/delete expected time

assuming update seq is indep of rand. choices

can we get good worst-case query time?

Assume static ...

Alternative analysis:

$$E \left[\text{total \# colliding pairs } (x,y) \in S \times S \right]$$

$$= E \left[\sum_{y \in S} (\# \text{ collisions with } y) \right]$$

$$= E \left[\sum_{y \in S} O\left(\frac{n}{m}\right) \right]$$

$$= O\left(n \cdot \frac{n}{m}\right).$$

$$\text{Set } m = cn^2 \Rightarrow E[\text{total \# colliding pairs}] = \underline{O\left(\frac{1}{c}\right)}$$

$$\Rightarrow \Pr[\text{total \# colliding pairs} \geq 1] \leq \underline{O\left(\frac{1}{c}\right)}.$$

by Markov's inequality

$$\Rightarrow \begin{array}{l} O(n^2) \text{ space} \\ O(1) \text{ worst-case query time} \end{array}$$

Perfect hashing

$$O(n) \text{ expected preproc time} \\ \text{by repeating } \underline{O(c)} \\ \text{expected \# times.}$$

Can space be reduced ^{back} to $O(n)$?

yes ... to be cont'd ...