

Lecture 3 (Advanced Algos) Jan 25, 2012

Hashing

Problem: Dictionary • preprocess set $S \subset U$, $|S|=n$.
↑
"universe"

• query: given x , report if " $x \in S$ "

Goal: fast query, small space.

$U =$ all IPs. 2^{32}

$S \subseteq U$.

Solutions:

1) store S
query: scan $S \rightarrow O(n)$ time

2) Binary Search Tree on S
space: $O(n)$
query: $O(\lg n)$.

3) store index: table $T[1..U]$

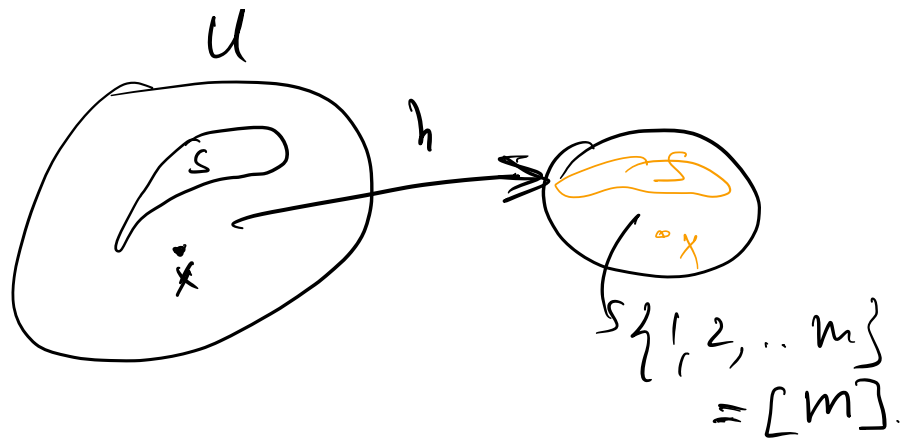
$T[i] = 1$ iff $i \in S$.

query: $O(1)$.

space: $O(U)$ bits.

Goal: obtain $O(1)$ q.t. $\leftarrow O(n)$ space

Hashing:



No Collision Prop: (NCP): $\forall x \in U, \forall y \in S, x \neq y:$
 $h(x) \neq h(y).$

Sol 4 (NCP): store table T s.t.

$$T[h(i)] = 1 \text{ iff } i \in S.$$

$$T[h(x)] = 1 \text{ iff } x \in S \text{ (NCP).}$$

Space: $\Theta(m).$

$$m = 10n.$$

Fact: $\exists h$ that satisf NCP. $m = n + 1.$

$$S = \{i_1, i_2, \dots, i_n\}$$

$$h(i_1) = 1$$

$$h(i_2) = 2.$$

\vdots

$$h(i_n) = n.$$

$$h(x) = n+1 \quad \forall x \notin S.$$

Issue: computing $h(x)$?

as hard as Dict. problem.

Sol 5: pick h indep. of S
at random.

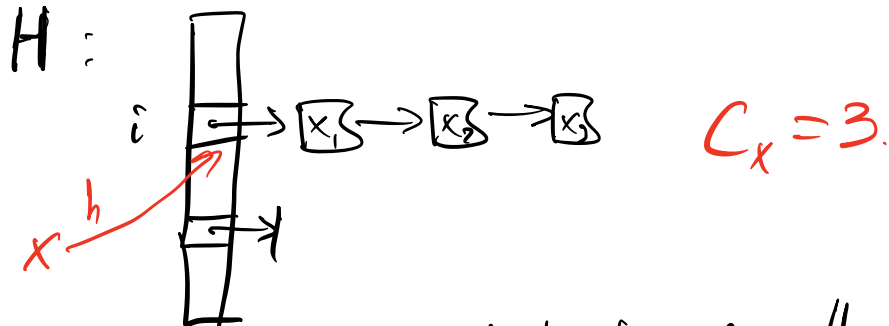
$\mathcal{H} = \{\text{all hash functions } h: U \rightarrow [m]\}$.

$$|\mathcal{H}| = m^{|U|}.$$

$h \in_r \mathcal{H}$.

Data structures: NCP does not hold
(with large prob.)

app: use chaining.



$H[i]$ = stores a linked-list of all
 y 's s.t. $h(y) = i$.

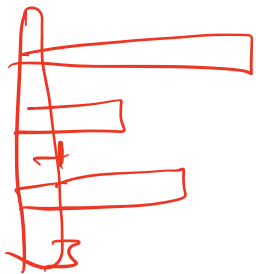
query: look-up L.L. of $h(x)$, iterate
and check if $\neq x$.

query time: prop to size of bucket $h(x)$.

$C_x \triangleq$ # elements y in S , different from x ,
s.t. $h(x) = h(y)$.

Q. 1. $= O(C_x)$.

fix x : $E_h[C_x] = E_h \left[\sum_{\substack{y \in S \\ y \neq x}} \mathbb{1}[h(x) = h(y)] \right]$
 $\left[\max_x E_h[C_x] \right]$
 indicator of event $h(x) = h(y)$.



$$= \sum_{\substack{y \in S \\ y \neq x}} \left[E_h \mathbb{1}[h(x) = h(y)] \right]$$

$$= \sum_{y \in S, y \neq x} P_r[h(x) = h(y)]$$

$$= \sum_{y \in S, y \neq x} \frac{1}{m} \leq \frac{n}{m}$$

Conclusion: expected q. time: $O(n/m + 1)$.

$\Rightarrow O(1)$ if $m = 2n$.

Remark: what is $E_h \left[\max_x C_x \right] = ?$

what is the size of largest bucket / linked list?

answer: $\Theta\left(\frac{\lg n}{\lg \lg n}\right)$ ← $m = \Theta(n)$.

Issue 1: how to choose/store $h \in \mathcal{H}$?

- If 2: worst-case runtime is \gg cont.
 \approx logarithmic.

Sol 6: (no more assumptions)

how to choose $h \in \mathcal{H}$. fully random.

Idea: use smaller / more structured families \mathcal{H} .

Def: \mathcal{H} is universal iff $\forall x \neq y$:

$$\Pr_{h \in \mathcal{H}} [h(x) = h(y)] = 1/m.$$

Remark: our analysis from Sol 5 used only fact that \mathcal{H} is universal.

Def: \mathcal{H} is k -wise independent iff.

\forall distinct $x_1, \dots, x_k \in U$, any. vals $v_1, \dots, v_k \in [m]$:

$$\Pr_{h \in \mathcal{H}} [h(x_i) = v_i, i = (1..k)] = \frac{1}{m^k}.$$

Fact: if \mathcal{H} is 2-wise indep \Rightarrow universal.

$$\forall x \neq y \quad \Pr [h(x) = h(y)] = \sum_v \Pr [h(x) = v \& h(y) = v] = m \cdot \frac{1}{m^2} = 1/m.$$

Fact: $\forall k$, there exist k -wise indep. h.f.

\mathcal{H} with $|\mathcal{H}| = O(k \cdot \lg U)$.

and time to compute $h(x)$ is \rightarrow

$$\text{2-wise: } h_{a,b}(x) = (a + bx) \% m.$$

Def: \mathcal{H} is α -universal, $\alpha \geq 1$, iff

$$\forall x \neq y: \Pr_h [h(x) = h(y)] \leq \frac{\alpha}{m}.$$

Remarks: $\mathbb{E}[C_x] \leq \alpha \cdot \frac{n}{m}$ ($h \in \mathcal{H}$ α -univ.)

Example of 2-universal \mathcal{H} :

$a \in [U]$ at random, odd.

$$h_a(x) = \lfloor ((a \cdot x) \% |U|) \cdot \frac{m}{|U|} \rfloor$$

$$h_a: U \rightarrow \{0, 1, \dots, m-1\}.$$

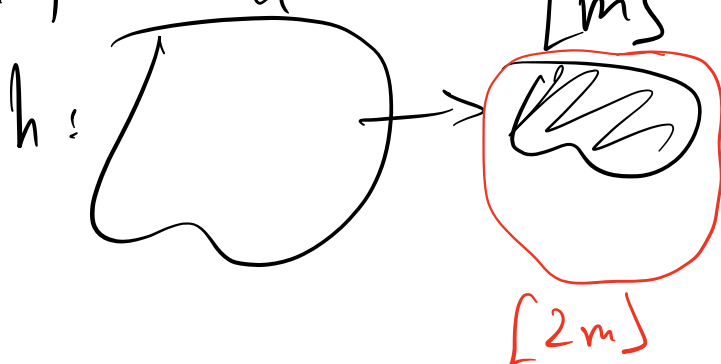
$$\mathcal{H} = \{h_a \text{ for } a \in [U] \text{ and odd}\}.$$

Conclusion (Sol 6): set $m=n$.

pick $h \in \mathcal{H}$ 2-univ.

space: $O(\lg U + n \cdot \lg U)$ bits to store
h.f. hash table hash table.

q.f: $O(1) + O(1)$ expected. (word operations)
h.f. eval. look-up. U $[m]$



Sol 7: Perfect Hashing.

Goal: get $O(1)$ q.f. deterministically.

idea: trade-off hash func construct
description
to depend (a little) on S

vs # collisions.