

Hashing 3

Feb 25 2021

Construction 3 (using finite fields)

$\text{GF}(2^u)$

each element

u -bit vector

Pick two random elements

$a, b \in \text{GF}(2^u)$

(uniformly at random)

$x \in U$

$$h(x) = ax + b$$

Q: 2-wise indep?

$$ax_1 + b = \alpha_1$$

$$ax_2 + b = \alpha_2$$

\Rightarrow

$$a =$$

$$b =$$

$$[u] \rightarrow [u]$$

we have $[U] \rightarrow [U] \rightarrow$ vector of length u

we want $[U] \rightarrow [M] \rightarrow$ vector of len. m ($m < u$)

$$|U| = 2^u \quad \uparrow \quad |M| = 2^m$$

$$|M| \ll |U| \\ m \ll u$$

Truncate to m bits! (Ex)

k -wise indep.

RF(2^u)

linear fn \rightarrow polynomial fn!

Pick uniformly at random: $a_0 \dots a_{k-1} \in \mathbb{R}_F(2^u)$

$$h(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{k-1} x^{k-1}$$

To prove: k -wise indep.

Other examples: Tabulation

Hash tables: closed addressing
Open addressing

Open addressing:

No separate DS

All keys stored in one array.

Linear probing:

$h(x)$

$g_0(h(x)+1) \bmod M$

Can have constant time look-ups.

but universal hashing not sufficient

5-wise indep is necessary

other probe sequences: step-size, quadratic etc.

Cuckoo Hashing;

Open addressing

by Pagh & Rodler (2004)

Take two tables T_1 & T_2

Both of size $M \approx O(N)$

Two hash fns $h_1, h_2: U \rightarrow [M]$

$h_1, h_2 \in H \leftarrow$ hash fn. family

Assume H is fully-random ($O(\log N)$ -wise indep)

Insertion of an element x :

$T_1[h_1(x)]$ or $T_2[h_2(x)]$

Bumping out process

\Rightarrow more than $6 * \log N$ bumps rehash everything.

Theorem: Expected time for insertion is $O(1)$ if $M \geq 4N$
table size \nearrow \nearrow (S)

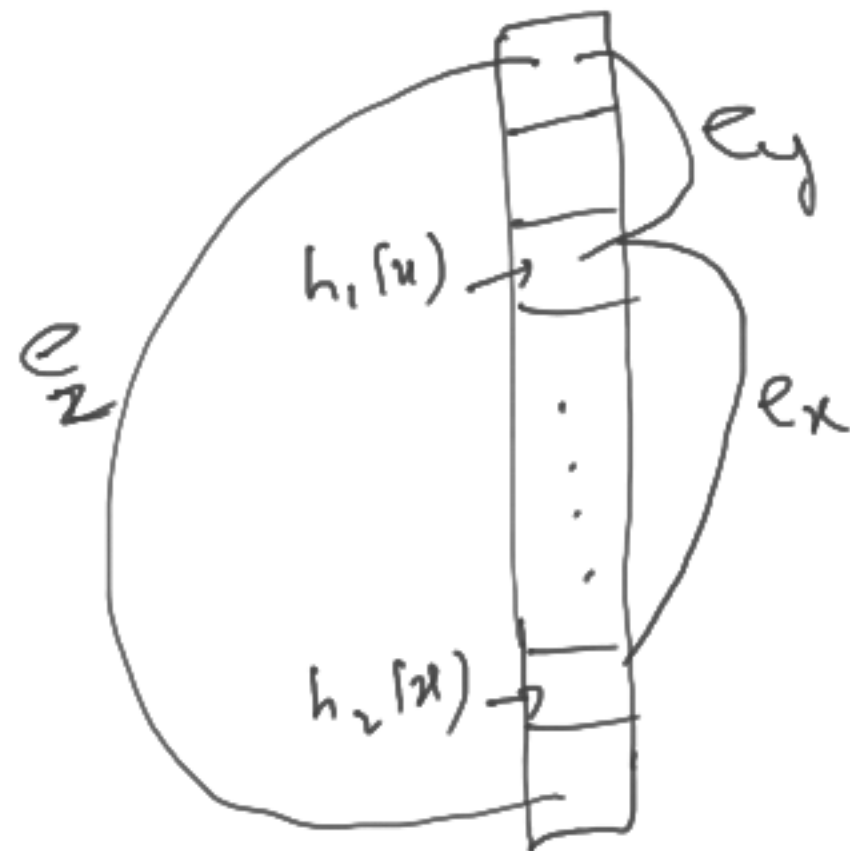
Proof sketch:

Cyclic graph (G) :

M vertices

Edges: elements to be inserted

$e_x = (h_1(x), h_2(x))$



Bucket of x

$B(x) =$ set of nodes of G reachable from $h_1(x)$ or $h_2(x)$

$$E[\text{insertion time for } x] = E[|B(x)|]$$

$\leftarrow G$ is random
- so is $B(x)$

Goal: $E[|B(x)|] \leq O(1)$

$$E[|B(x)|] = \sum_{\substack{y \in S \\ y \neq x}} P[e_y \in B(x)]$$

$$\leq N P[e_y \in B(x)]$$

Sufficient to show $P[e_y \in B(x)] \leq O(1/M)$

Lemma: For any i, j in $[M]$

$$P[\text{there exists a path of length } l \text{ between } i \text{ \& } j \text{ in the circled graph } \mathcal{G}_l] \leq \frac{1}{2^l M}$$

Proof: Via induction.

$$l=1: P(\text{edge between } i \text{ \& } j)$$

$$= P(\exists y \text{ s.t. } e_y \text{ exists in } \mathcal{G}_1)$$

$$\leq N \cdot \frac{2}{M^2}$$

$$\leftarrow P\left[\begin{aligned} & (h_1(y)=i \cap h_2(y)=j) \\ & \cup h_1(y)=j \cap h_2(y)=i \end{aligned}\right]$$

$$\leq \frac{1}{2 \cdot M}$$

Then induction on l .
(Ex)

To show: $P[e_y \in B(x)] \leq O(1/M)$

using the lemma,

$$P[e_y \in B(x)] \leq \sum_{l \geq 1} \frac{1}{2^l M}$$
$$= O(1/M).$$

$M \geq 4N \Rightarrow 25\%$ space efficiency

For $d=3$, experimentally $>90\%$ space efficiency

Bloom Filters:

Limited ops; membership query

Allows for mistakes:

only false positives; no false negatives.

• Array T of M bits
- initialized to 0

• k hash func. $h_1, h_2, \dots, h_k: V \rightarrow [M]$
- assume fully random

• Adding: $x \in S$
set bits $T[h_1(x)] \ T[h_2(x)] \ \dots \ T[h_k(x)]$ to 1

Membership query: