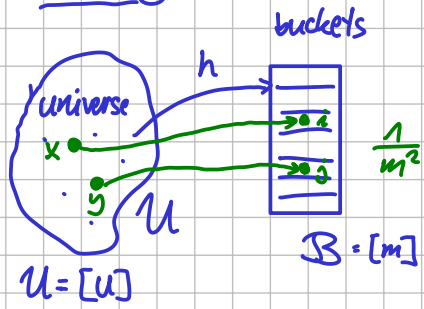


Hashing



choose $h \in \mathcal{H}$ randomly
For hashing with buckets

$E[\text{time complexity}] \in O\left(\frac{n}{m}\right)$
↑
 $h \in \mathcal{H}$

Df: A family \mathcal{H} of functions from U to B is c -universal for $c > 0$ iff

$$\forall x, y \in U, x \neq y: \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq \frac{c}{m}$$

Df: A family \mathcal{H} of functions from U to B is (k, c) -independent for $k \in \mathbb{N}^+$, $c > 0$, $|U| \geq k$ iff

$$\forall x_1, \dots, x_k \in U \text{ all distinct } \forall b_1, \dots, b_k \in B$$

$$\Pr_{h \in \mathcal{H}} [h(x_1) = b_1 \& \dots \& h(x_k) = b_k] \leq \frac{c}{m^k}$$

Df: \mathcal{H} is k -independent iff $\exists c > 0: \mathcal{H}$ is (k, c) -indep.

- ☺ \mathcal{H} is (k, c) -ind. $\Rightarrow (k-1, c)$ -indep
- \mathcal{H} is $(2, c)$ -ind. $\Rightarrow c$ -universal
- \mathcal{H} is $(1, c)$ -ind.

Consider: $\mathcal{C} := \{h_a \mid a \in B\}$
 $h_a(x) := a$ for all x
Then \mathcal{C} is $(1, 1)$ -indep.

Df: For p prime and $m \leq p$:

$$\mathcal{L} := \{h_{a,b} \mid a, b \in [p]\}$$

linear functions

$$h_{a,b}(x) := ((ax + b) \bmod p) \bmod m$$

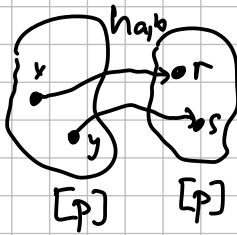
calculating in a field \mathbb{Z}_p

Df: $\mathcal{L}' := \{h_{a,b} \mid a, b \in [p], a \neq 0\}$

Theorem: \mathcal{L} is 2-universal.
 \mathcal{L}' is 1-universal

Proof: Let $x, y \in [p], x \neq y$.

First without mod m :



- params $(a, b) \in [p]^2$
 - $r := (ax + b) \bmod p$
 $s := (ay + b) \bmod p$ $r \neq s$
 - $(a, b) \xleftrightarrow{\pm 1} (r, s)$ pick (a, b) u.a.r.
pick (r, s) u.a.r.
- $p(p-1)$ pairs $p(p-1)$ pairs

With mod m :

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

$$= \Pr_{a,b} [((ax + b) \bmod p) \bmod m = ((ay + b) \bmod p) \bmod m]$$

$$= \Pr_{r,s} [r \equiv s \pmod{m}] \leq \frac{2}{m}$$

Df: (r, s) is a bad pair $\equiv r \equiv s \pmod{m}$

Count bad pairs -- for a given r



Sum over all r

$$\#s: (r, s) \text{ is BP} \leq \left\lceil \frac{p}{m} \right\rceil \leq \frac{p+m-1}{m} \leq \frac{2p}{m}$$

$$\#BP \leq \frac{2p^2}{m}$$

$$\#s \leq \left\lceil \frac{p}{m} \right\rceil - 1$$

$$\leq \frac{p+m-1}{m} - 1$$

$$= \frac{p-1}{m}$$

$$\# \text{ all pairs} = p^2 \Rightarrow \Pr[\text{a pair is bad}] = \frac{\#BP}{p^2} \leq \frac{2}{m}$$

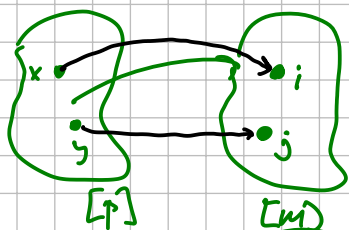
$$\#BP \leq \frac{p \cdot (p-1)}{m}$$

$$\Pr[\text{pair is bad}] \leq \frac{p \cdot (p-1)}{m} / p(p-1) = \frac{1}{m}$$

Theorem: \mathcal{L} is $(2,1)$ -independent.

we want this to be $\leq \frac{4}{m^2}$

Proof:



$$\Pr_{r,s} [r \equiv i \pmod{m} \ \& \ s \equiv j \pmod{m}]$$

$$= \underbrace{\Pr_{r,s} [r \equiv i \pmod{m}]}_{\leq 2/m} \cdot \underbrace{\Pr_{r,s} [s \equiv j \pmod{m}]}_{\leq 2/m} \leq \frac{4}{m^2}$$

independent events

in each of $\lceil \frac{p}{m} \rceil$ intervals, there is at most one such r

$$\frac{\lceil \frac{p}{m} \rceil}{p} \leq \frac{p+m-1}{m \cdot p} \leq \frac{2p}{mp} = \frac{2}{m}$$

Construction of k -indep. families for arbitrary k :

Df: For $U = [p]$ and $k \geq 1$:

$\mathcal{P}_k := \{h_{\vec{a}} \mid \vec{a} \in [p]^k\}$ family of functions from $[p]$ to $[p]$

$$h_{\vec{a}}(x) := \sum_{i=0}^{k-1} a_i \cdot x^i \pmod{p}$$

← all this mod p

reduce mod m using general lemma later

Theorem: \mathcal{P}_k is $(k,1)$ -independent

proof sketch:

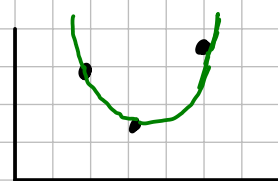
$$\forall i \ h_{\vec{a}}(x_i) = b_i$$

unique poly of deg $< k$

unique choice of \vec{a}

$$\Pr[\forall i \ h_{\vec{a}}(x_i) = b_i] = \frac{1}{p^k}$$

Given $x_1 \dots x_k$ items and $b_1 \dots b_k$ buckets



for k points in the plane $\exists!$ polynomial of deg $< k$ connecting them

Good news: this works for arbitrary k

Bad news: time is $\Theta(k)$

Multiply-shift hashing

$$a \cdot x$$

↑ odd

32-bit machine



Df: $x \langle i:j \rangle ::$ bits i to $j-1$ of x

$$a * x \gg (32-l)$$

$$(a * x) \langle 32-l:32 \rangle$$

top most l bits



Df: For given w (word size) and l (result size) in bits

$\mathcal{M} := \{h_a \mid a \in [2^w], a \text{ odd}\}$ from $[2^w] \rightarrow [2^l]$

$$x \langle 0:1 \rangle = \text{least-sig. bit} = x \bmod 2$$

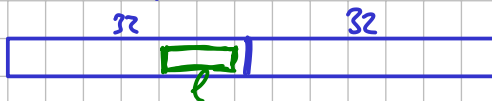
$$x \langle 0:8 \rangle = \text{lowest 8 bits..}$$

$$h_a(x) := (ax) \langle w-l:w \rangle$$

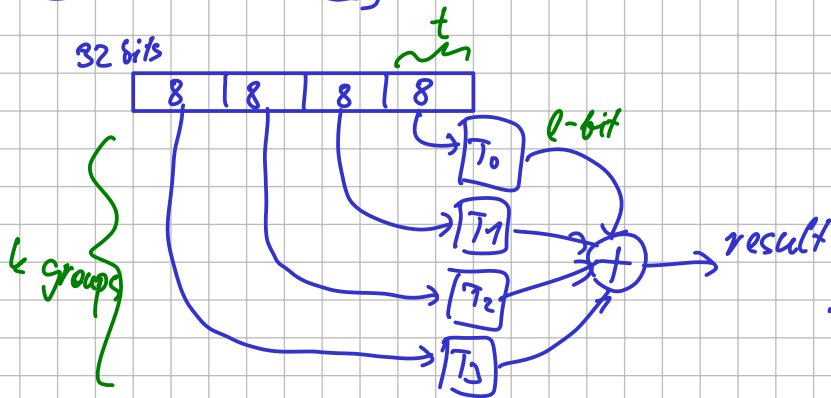
$$x \langle i:j \rangle = \lfloor x/2^i \rfloor \bmod 2^{j-i}$$

Claim: \mathcal{M} is 2-universal.

$\rightarrow \exists$ version of mult-shift which is 2-indep.



Tabulation Hashing



k groups per l bits

$$[2^{kt}] \rightarrow [2^l]$$

tables T_0, \dots, T_{k-1}

each $T_i: [2^l] \rightarrow [2^l]$

Then

$$h_{T_0, \dots, T_{k-1}}(x) := \bigoplus_{0 \leq i < k} T_i[x \ll (i+1)l]$$

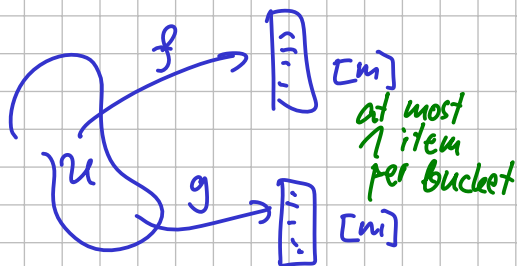
Time for evaluating h is $\Theta(k)$.

Memory needed: $\Theta(k \cdot 2^l \cdot l)$

Claim: 3-independent, but not 4-independent.

Cuckoo Hashing

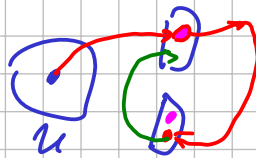
① 2 tables, 2 functions



x is stored either at $B_1[f(x)]$ or at $B_2[g(x)]$

Find is $O(1)$ in worst case!

Insert:



keep kicking out items & moving to the other location until

the location is empty

timeout

rehash with new f, g

② 1 table, 2 functions



Theorem: Let $\epsilon > 0$ be a fixed constant,

$$m \geq (2 + \epsilon) \cdot n,$$

insertion timeout $\lceil 6 \cdot \log m \rceil$

f, g uniformly at random from $\lceil 6 \log m \rceil$ -indep. family.

Then $\mathbb{E}_{f, g}[\text{time for insert}] \in O(1)$.

known:

6-indep. is insufficient but tabulation works!

depends on ϵ