

Lecture 8:

①

Last time: A set \mathcal{H} of hash functions $\{0, 1, \dots, U-1\} \rightarrow \{0, 1, \dots, n\}$ is universal if, \forall distinct $x, y \in U$,

$$\Pr_{h \in \mathcal{H}} [h(x) \neq h(y)] \leq \frac{1}{n}. \quad [\leftarrow \text{always } = \frac{1}{n} \text{ if } \mathcal{H} \text{ is "all functions"}, \text{ as in SUHA}]$$

(Comment: "c-almost-universal" has " $\leq \frac{C}{n}$ " in place of " $\leq \frac{1}{n}$ ".

If $c=2$, e.g., morally just as good & often

easier to achieve.]

[We saw an example last time showing that universal hash families can be as good as SUHA for at least some nice hashing apps. We'll see more today, but first, let's construct them...]

example 1: [Not the most efficient, but simple proof.]

$$\text{Assume } U = 2^u \quad (\text{e.g.: } u=64, 128)$$

$$n = 2^l \quad (\text{e.g.: } l=16 \quad (\text{table of size 65536}))$$

~~This construction~~ Suppose $A = \begin{bmatrix} 000 & \dots & \\ 011 & \dots & \\ \vdots & \ddots & \vdots \\ & \dots & \end{bmatrix}$ is a matrix of bits. [couple thousand bits].

$$\text{Define } h_A(x) = \begin{bmatrix} \text{l-bit string} \\ A \\ x \end{bmatrix} = \begin{bmatrix} \text{l-bit string} \\ \mod 2 \\ \text{(in each entry)} \end{bmatrix}$$

\downarrow l-bit string \equiv hash slot in $0 \dots n-1$

~~obviously~~, i^{th} ^{bit} ~~entry~~ of $h_A(x)$ is $a^{(i)} \cdot x \bmod 2$ $\left[\begin{array}{l} \text{computable in a few machine instructions, esp. if } a^{(i)} \\ \text{fits in a word} \end{array} \right]$

$\mathcal{H} = \{h_A : A \in \{0, 1\}^{l \times u}\}$. Picking $h \in \mathcal{H} \equiv$ picking A . $\left[\begin{array}{l} h_A(x) \text{ computable in } O(l) \end{array} \right]$

Claim: This is universal.

②

Proof: Fix any $x, y \in \{0,1\}^l$. Need to show $\Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{n}$

$$\Leftrightarrow \Pr_A [Ax = Ay \pmod{2}] \leq \frac{1}{n}$$

$$\Leftrightarrow \Pr_A [A(x-y) = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \pmod{2}] \leq \frac{1}{n}.$$

$$\Leftrightarrow \Pr_A [Az = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \pmod{2}] \leq \left(\frac{1}{2}\right)^l$$

where $z = x-y \pmod{2}$
 $\neq \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \because x \neq y$.

$\Rightarrow z$ has at least one 1, say z_i .

$$\begin{bmatrix} | & | & & | \\ a_1 & a_2 & \cdots & a_u \\ | & | & & | \\ A & & & \end{bmatrix} \begin{bmatrix} * \\ * \\ \vdots \\ 1 \\ + \\ * \\ \vdots \\ * \end{bmatrix} \stackrel{?}{=} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \pmod{2}.$$

$$Az = z_1 \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_i \\ \vdots \\ a_u \end{bmatrix} + z_2 \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_i \\ \vdots \\ a_u \end{bmatrix} + \dots + z_u \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_i \\ \vdots \\ a_u \end{bmatrix} \quad \text{where } \begin{bmatrix} a_j \end{bmatrix} \text{ is all random.}$$

Imagine all a_j 's except a_i picked first, then a_i picked.

$$Az = \cancel{\begin{bmatrix} | & | & & | \\ a_1 & a_2 & \cdots & a_u \\ | & | & & | \\ A & & & \end{bmatrix}} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \begin{bmatrix} | \\ a_i \\ | \end{bmatrix} \stackrel{?}{=} \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \pmod{2} \quad \text{Probability} = \left(\frac{1}{2}\right)^l \text{ exactly over remaining choice of } a_i$$

✗

Con 1: Storing A : lu bits. [Okay, but $2u$ is possible]

Con 2: Evaluating $h_A(x) = \Theta(l)$ word operations [

[O(1) is possible]

Con 3: $h_A(0) = 0$ always!

[Kinda weird! "Each 1 item goes to rand. place" fails.]

Cheap fix for Cons 3: Also pick random $b \in \{0,1\}^l$, ③

let $h_{A,b}(x) = Ax + b \bmod 2$

Still universal ✓ $\quad [h_{A,b}(x) = h_{A,b}(y) \Leftrightarrow Ax + b = Ay + b \bmod 2]$
 $\qquad \qquad \qquad \qquad \qquad \Leftrightarrow Ax = Ay \bmod 2,$
so same proof. ✓

"Uniformity": For all $x < 2^l$, for all $j < n$, $\Pr_{h=h_{A,b}}[h(x)=j] = \frac{1}{n}$.

Because: $h_{A,b}(x)=j \Leftrightarrow \begin{bmatrix} A \\ b \end{bmatrix} \left| \begin{bmatrix} x \\ j \end{bmatrix} \right. + \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} j \\ 1 \end{bmatrix} \bmod 2$

Pick A first, b last. $\Rightarrow \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} b \\ 0 \end{bmatrix} = \begin{bmatrix} j \\ 1 \end{bmatrix} \bmod 2$
 $\begin{bmatrix} 0 \\ 1 \end{bmatrix} \xrightarrow{Ax} \begin{bmatrix} b \\ 0 \end{bmatrix}$
Prob is $\left(\frac{1}{2}\right)^l = \frac{1}{n}$. ✓

Fix for Cons 1 & 2:

Carter & Wegman's method:

- ~~Fix~~ any prime number p between 2^l & 2^{2l} // takes a bit of effort, but it's a one-time cost.
Fix // doesn't have to be random //

To choose h : choose $a \in \{1, 2, \dots, p-1\}$ // $\approx l$ bits
 $b \in \{0, 1, \dots, p-1\}$ // l bits
at random.

$$h_{a,b}(x) := ((a \cdot x + b) \bmod p) \bmod n.$$

// Note: can't skip the mod n part! exercise

exercise:

{ kinda annoying,
but elementary //

This is universal.

Pro 1: Storing a, b : $\approx 2l$ bits

Pro 2: Eval. $h_{a,b}(x)$: $O(1)$ (?) machine ops.

More options:

- Pick odd $a \in \{1, 3, 5, 7, \dots, n-1\}$ at random
- $h_a(x) := \underbrace{(a \cdot x \bmod u)}_{\substack{\text{drop all but } u \\ \text{LSBs}}} \quad \underbrace{\text{div } \frac{u}{n}}_{\substack{\text{right-shift by } u-l \text{ bits}}}$

Pro 1: Storing a : l bits \Leftarrow

Pro 2: Eval $h_a(x)$: $O(1)$ ops (3?)

Con: Not universal. But is "2-almost-universal." $\cancel{\text{X}}$

Proof: exercise
[not too hard]

Universal + Uniform: For any two items $x \neq y$,

the ~~2~~² bins ~~are~~ $h(x), h(y)$
are a completely random pair.

[each possible
occurs with
prob $\frac{1}{n^2}$]

Higher analogue: " k^3 -wise ~~independent~~ hash family"

For any ^{distinct}₃ items x_1, x_2, x_3 , $(h(x_1), h(x_2), h(x_3))$
is each possibility with prob. $\frac{1}{n^3}$.

Fact: • \exists efficient-to-eval k -wise indep. hash families

where you store $O(k \cdot u)$ bits. [but a bit elaborate to explain $\cancel{\text{X}}$]

• 5-wise indep. is good enough for many SUHA-achievable properties

$\cancel{\text{X}}$

One more example where universal (or 2-almost-universal) is "good enough": "Perfect Hashing".

Goal: Worst-case $O(1)$ lookups, space $O(m)$...

5

[[sounds awesome, better than w/ SUHA. Catch is...]]
... for a static hash table:

All m items $x_1 \dots x_m$ given in advance, no ~~insert/dels,~~ just lookups.
[[Could build an optimal BST or similar, but hashing will suffice.]]

Say we use a universal family with $n=m$. [[Optimistic; last time we

Recall: After hashing, if $C_{ij} = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \neq \end{cases}$, considered $n \gg m^2$]]

$$C = \sum_{1 \leq i < j \leq n} C_{ij} = \# \text{ colliding pairs},$$

$$\mathbb{E}[C_{ij}] \leq \frac{1}{n} \quad (\text{universality}), \quad \mathbb{E}[C] \leq \binom{m}{2} \cdot \frac{1}{n} \quad [\text{in. of espec.}]$$
$$\leq \frac{m^2}{2n} = \frac{m}{2} \quad (\text{for } n=m)$$

Q: If 5th bin has load 4, what will C be (at least?)

$\sum_{s=0}^{4-1} \binom{4}{s}$

Markov: $\Pr[C \geq 20 \cdot \frac{m}{2}] \leq \frac{1}{20}; \therefore C \leq 10m \text{ except w. prob.} \leq 5\%$

Is that... good? Bad? Last time we said $C=0 \Rightarrow$ no collisions
What could $C \leq 10m$ imply?

Say we get...

$$\begin{array}{cccc} \underbrace{\begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix}}_{L_0} & \underbrace{\begin{matrix} 0 \\ 0 \end{matrix}}_{\substack{(2) \\ \text{colliding} \\ \text{pairs}}} & \underbrace{\begin{matrix} 0 \\ 0 \end{matrix}}_{\substack{(2) \\ \text{colliding} \\ \text{pairs}}} & \underbrace{\begin{matrix} 0 \\ 0 \end{matrix}}_{\substack{(2) \\ \text{colliding} \\ \text{pairs}}} \\ L_1 & L_2 & L_3 & L_4 \\ \binom{5}{2} = 10 & \binom{3}{2} = 3 & \binom{4}{2} = 6 & \binom{1}{2} = 0 \end{array} \dots$$

What's C ?

$$\therefore C = \sum_{i=0}^{n-1} \binom{L_i}{2}, \quad L_i := \text{load of } i^{\text{th}} \text{ bin.}$$

Note: $\binom{L_i}{2} = L_i \frac{(L_i-1)}{2} \approx \frac{L_i^2}{2}$. \therefore If some $L_i > 5\sqrt{m}$, then $L_i^2 > \frac{25m}{2} \gg 10m$
 $\Rightarrow C > 10m$, contrary to *

\therefore except with prob $\leq 5\%$, Max load $\leq 5\sqrt{m}$.

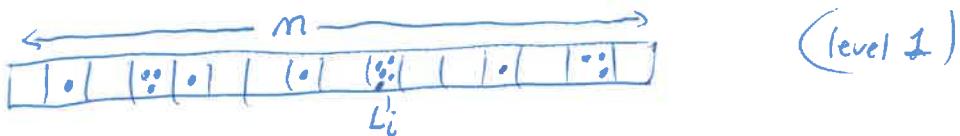
Not totally awesome. For $m=1$, w.h.p. the max load is $O(\frac{\log m}{\log \log m})$ under SHTA. STM is way bigger.

But... Moreover, except with prob $\leq 5\%$,

$$10m \geq C = \sum_{i=0}^{n-1} \binom{L_i}{2} \geq \sum_{i:L_i \geq 2} \frac{L_i^2}{4}, \quad \therefore \binom{L_i}{2} \geq \frac{L_i^2}{4} \text{ when } L_i \geq 2$$

$$\Rightarrow \sum_{i:L_i \geq 2} L_i^2 \leq 40m.$$

Idea: Hash with universal family, compute L_i 's. Check that ~~**~~ holds
IF not, pick a new h and try again. But happens 95% of time?



Now for each slot i with $L_i \geq 2$,



build a second-level hash table of size $10L_i^2$ with a new random hash function h_i from a universal family

Last time we saw: for universal hashing,

if L_i items $\rightarrow 10L_i^2$ slots, then $\geq 95\%$ of the time, max load is 1.

So do this for each level-1 slot with load ≥ 1 .

Again, after hashing the L_i items, check load is ≤ 1 . If not, rehash, but $\leq 5\%$ chance.

Space: m (level 1)

$$+ \sum_{i:L_i \geq 2} 10L_i^2 \leq 400m, \text{ by } \del{**}$$

\therefore total space: $O(m)$

Worst-case lookup: 2 hashes.

(Slight bummer: needs up to m different hash functions
 $\rightarrow m \cdot O(1)$ bits to store these.)