# Lecture 4    1/27/22.

## Perfect Hashing

Goal: Dict. with $O(n)$ space

$O(1)$ query time deterministic.

Contrast: last time $O(1)$ expected.

Question: what if build a hash table
with random h.f with no collisions.

$$C_x = \# y \in S \quad y \neq x \text{ and } h(y) = h(x).$$

$$C \overset{\Delta}{=} \sum_{x \in S} C_x = \# \text{ pairwise col. } x, y \in S.$$

Suppose we want $C = 0$. $\Rightarrow$ all buckets
have size 0 or 1.

$E[C]$ when we have $m$ buckets.

$$\boxed{h: U \to [m]}.$$

SCU
$n = |S|$
$m = $ table size

$$E[C] = E\left[\sum_{x \in S} C_x\right] = \sum_{x \in S} E[C_x]$$

$$\leq \sum_{x \in S} \frac{n}{m}$$

$$= n^2/m.$$

By Markov bnd: $C \leq 4 \mathbb{E}[C]$ with prob $\geq 1 - \frac{1}{4}$.

$$= 4n^2/m.$$

Set $m = 8n^2 \Rightarrow$ $C \leq \frac{1}{2}$ with prob $\geq 1 - \frac{1}{4}$.

$$\Rightarrow C = 0 \quad \text{———} \quad \text{//———}$$

Algo: repeat until success:
- pick random $h \in \mathcal{H}$, $m = 8n^2$
- compute $C$
- if $C \geq 1$, repeat.
- if $C = 0$, success.

build hash table $H$ with last $h$.

$$\text{space } O(m+n) = O(n^2)$$

Algo has space $O(m+n) = O(n^2)$

q.l. $O(1)$ determ.



Preprocessing: proportional to #repeats:
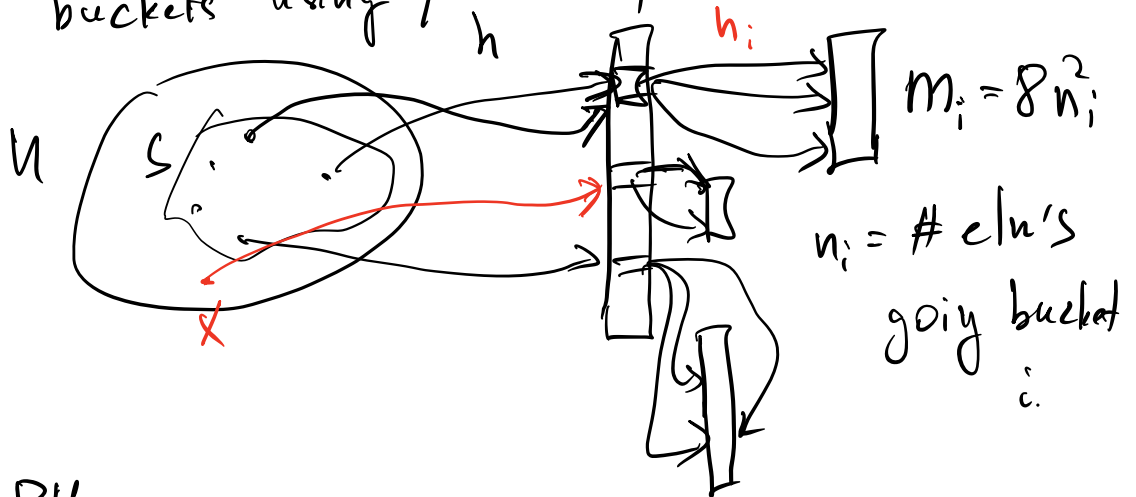
$$\mathbb{E}[\#repeats] \leq 1 \cdot \frac{3}{4} + 2 \cdot \frac{1}{4} \cdot \frac{3}{4} + 3 \cdot \left(\frac{1}{4}\right)^2 \cdot \frac{3}{4} + \ldots$$

$$= 1 + \frac{1}{4} + \left(\frac{1}{4}\right)^2 + \ldots$$

$$= \frac{1}{1 - \frac{1}{4}} = 4/3.$$

Con: can obtain $O(n^2)$ space, $O(1)$ det. q.t.

Perfect Hashing

<u>Idea:</u> use standard hashing, $m=2n$.
$+ 2^{nd}$ level hash where we split non-temp
buckets using quadratic-space hash.



$m_i = 8n_i^2$

$n_i = $ # eln's
going bucket
$i$.

<u>Algo PH:</u>

Pick a random h.f. $h: U \to [m]$, $m=2n$.

for $i=1..m:$
- build $2^{nd}$ level hash table of size
  $m_i = 8n_i^2$, (using <u>Cor</u>)
  where $n_i = |S \cap h^{-1}(i)| = $ # eln's in
  bucket $i$.
- store $S \cap h^{-1}(i)$ in the $2^{nd}$ level h.t.

<u>Query time:</u> look-up $h(x)$, then look-up
hash table corresponding to $i=h(x)$.

<u>Obs 1:</u> query time is $O(1)$ det.

Question: what about space?

<u>Claim:</u> in expectation, size of $2^{nd}$ level hash tables is $O(n)$.

$n_i + n_i(n_i - 1) = n_i^2$

<u>pf:</u> $\sum\limits_{i=1}^{m} m_i = \sum\limits_{i=1}^{m} 8n_i^2$.

$$\mathbb{E}_h \left[ \sum\limits_{i=1}^{n} n_i^2 \right] = \mathbb{E}_h \left[ \sum\limits_{\substack{i=1 \\ n_i \geq 1}}^{m} n_i + n_i(n_i - 1) \right]$$

$$= \mathbb{E} \sum\limits_{i=1}^{m} n_i + \mathbb{E} \left[ \sum\limits_{\substack{i=1 \\ n_i \geq 1}}^{m} n_i(n_i - 1) \right]$$

$$= n + \mathbb{E}[C]$$

$$= n + n^2/m = n + n/2 = O(n).$$

<u>Conclusion:</u> PH algo obtains $O(1)$ det q.t.

$O(n)$ space <u>exp.</u>

$O(n)$ preproc. <u>exp.</u>

<u>Remark:</u> can have $O(n)$ det. space bound, by repeating $1^{st}$ lev hashing until $\sum n_i^2 \leq O(n)$.
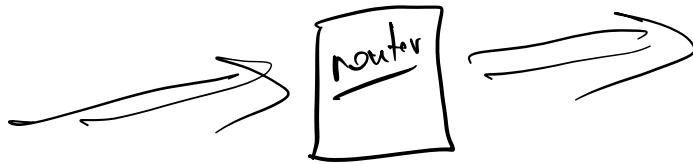
<u>Remark:</u> <u>dynamic Dict:</u> can we obtain $O(1)$ det. query/update time?

(with $O(n)$ space).  **OPEN.**

## Streaming & Sketching

Mod 1:



Mod 2: more eff.
read data in
a linear fashion.

Goal: store statistic on data, as streams
to be used later.

Distinct elm count.

Stream: elements from universe $[n]$
$$= \{1, .. n\}.$$

Eg.: $[n] =$ all possible IPs.

Problem: report how many diff elm's
we have seen in the stream.

Stream length: $m$.

Goal: store as little info as possible

Solution 1: store entire stream: $O(m)$ space.

Sol 2: store a table $T[1..n]$,

$T[i] = 1$ iff we've seen $i$.

$O(n)$ bits.

Can we obtain space $\ll \min\{n, m\}$?

No unless we allow approx + random.

**Algo: approx & random. [Flajolet-Martin]**

- pick a random hash func. $h: [n] \to [0,1]$.
- store a reg. $z$, init $z = 1$.
- when see $i$ in the stream:
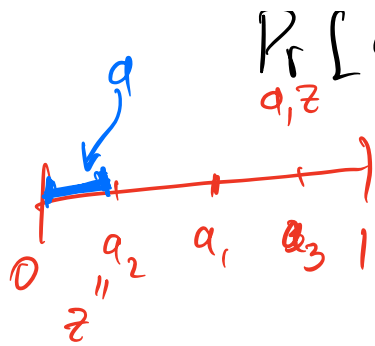$$z = \min\{z, h(i)\}.$$
- Est: $\frac{1}{z} - 1$.

**Note:** $z = \min$ hash function value $h(i)$ among $i$'s in the stream.

**Claim:** $\mathbb{E}\{z\} = \frac{1}{d+1}$, $d = \#$ distinct elm's.

**Pf:** $z = \min$ of $d$ random #'s $c_1 \dots c_d \in [0,1]$.
$$\underset{\parallel}{\phantom{c_1}} $$
$$h(1^{st} elm)$$

Experiment: pick $a \in [0,1]$.

$\nearrow$ 1) $z$ $\quad$ 1) $= \mathbb{E}\left[\Pr_a[a < z]\right]$

$$Pr_{a,z}[a < z]$$



$\to$ 2) Prob that $a$ is min $a_1, a_2, \dots a_d, a$

$$= \frac{1}{d+1}$$

$$\Rightarrow E[z] = \frac{1}{d+1} \quad \boxtimes.$$

$$E\left[\frac{1}{z}\right] \neq \frac{1}{E[z]}.$$