

# Lecture #4 (AA). 1/21/21.

## Perfect Hashing

Recap: Dict: given  $S \subseteq U$ , query "is  $x \in S$ ".

$$n \triangleq |S|$$

$h: U \rightarrow [m]$ . random / pseudo-random (can

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

$$C = \sum_{x \in S} C_x.$$

Claim (last time):  $E[C] \leq n^2/m$ .

Hence can construct a hash table of size

$m = 4n^2$  where each bucket size  $\leq 1$ .

larger than before.

query time =  $O(1)$   
always.

Thm [P.M.]: can get d.s. space  $O(n)$ ,

q.t. =  $O(1)$  always. Construction time

=  $O(n)$  in  
expected time

Pf: just combine "standard" hashing

+  $4n^2$  solution.



Issue: some buckets may be large  $\gg 1$ .

$n_i \triangleq$  size of bucket  $i = |h^{-1}(i) \cap S|$ .

Add 2<sup>nd</sup> level using P.H. quadratic-SP sol.

Algo:

- pick random  $h: U \rightarrow [m]$ ,  $m \approx n$ .
- $n_i$  def as above.
- for each  $i \in [m]$ ,
  - build a 2<sup>nd</sup> level h.t., using  $h_i: U \rightarrow [4n_i^2]$
  - such that no collisions in  $2^{n_i}$

also from past lecture

Analysis. 1. Correctness immediate.

2. Query time:  $O(1)$  always by construction.

3. Space =  $O(n)$  +  $O(\sum_i n_i)$ .

1st lev. h.t.

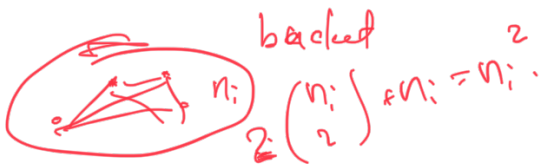
2nd level h.t.'s.

Claim:  $E[\sum_i n_i^2] = O(n)$ .

Pf:  $E[\sum_i n_i^2] \leq E[n_i + \sum_{j=1}^{2^x} \text{collisions in bucket } j]$

$= E[n + \sum_i n_i] = n + E[\sum_i n_i]$

$n_i^2 = n_i(n_i - 1) + n_i = n_i + 2 \binom{n_i}{2}$

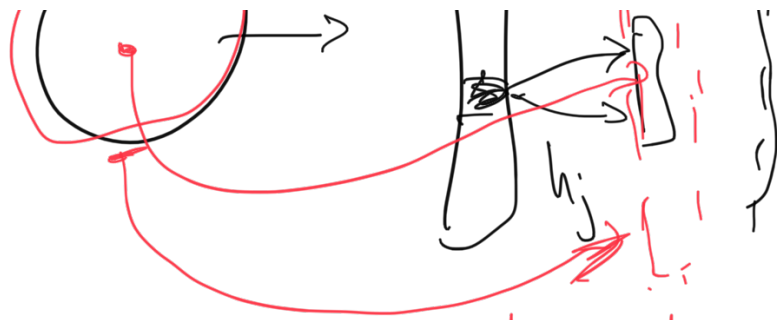


$= n + \frac{n^2}{m} = 2n$   $\square$

rad thm.

Obs: space is  $O(n)$  in expectation only. by Markov,  $O(n)$  with prob  $\geq 90\%$   $\Rightarrow$  can repeat construction  $O(1)$  exp. # times until space is  $O(n)$ .



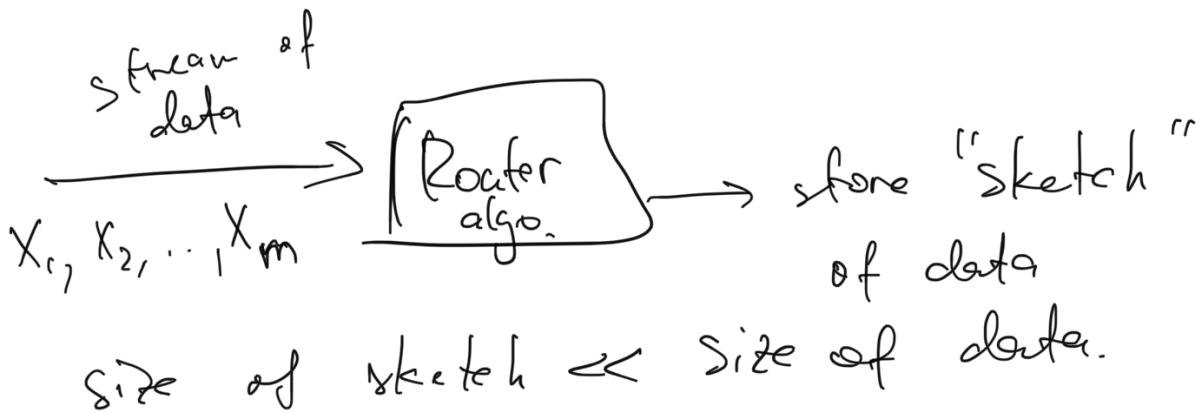


$h'$  depends on  $S$ .

Remark : - in all this construction, used  $\mathbb{F}$  to analyze.  
 $\Rightarrow$  enough to use universal hash. f. (both  $1^{st}$  &  $2^{nd}$ )


## II Sketching & Streaming algos.

Model: streaming model.



Few problems to consider.  
Problem: counting # distinct elements.  
 .. n. ?

$$X_1, \dots, X_m \in [n].$$

Performance measure (complexity) : space used by  to be able to answer query @ each step

1.  $m \cdot \lg n$  space (store every thing).  
 $O(d \cdot \lg n)$ . just store the distinct items.

2. via hashing:

$O(n)$  → bit-table

$O(d)$  → hash table of items seen so far.

$d = \#$  distinct items.

Thm: for  $\epsilon < 1/2$ ,  $d > \Omega(1/\epsilon^2)$ .

can get space  $O\left(\frac{\lg n}{\epsilon^2}\right)$  space

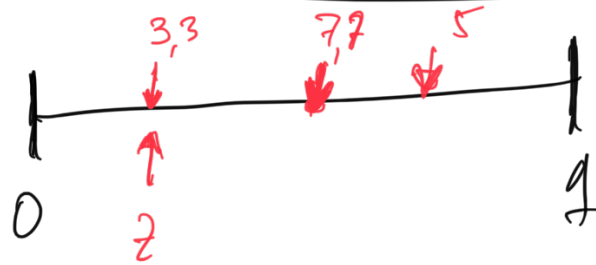
Algo [Flajolet-Martin '85]:

uses a hash function  $h: [n] \rightarrow [0, 1]$ .  
 fully random.

→ Init:  $z = 1$ .

- @ item  $i$ :  $z := \min\{z, h(i)\}$ .

- Estimator:  $\frac{1}{z} - 1$ .

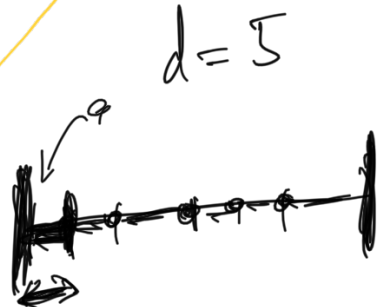


Obs  $z$  depends only on  $d$ .

$$E\left[\frac{1}{z} - 1\right] = +\infty.$$

Claim:  $E[z] = \frac{1}{d+1}$ .

pf:  $z = \min$  of  $d$  random vars  $\in [0, 1]$ .



Pick  $a \in [0, 1]$  at random.

$$\Pr\{a < z\} = \Pr\left[\text{choose } d+1 \text{ r.v. } \in [0, 1] \text{ last chosen one is smallest}\right]$$

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

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

Q

In general  $\mathbb{E}[\frac{1}{z}] \neq \frac{1}{\mathbb{E}[z]}$ .

Claim:  $\text{Var}[z] \leq z/d^2$ .

pf. omitted.

By Chebyshev:  $z = \frac{1}{d \epsilon} \pm O\left(\frac{1}{d}\right)$  with  $\approx 90\% p$

Possible to get  $1 \pm \epsilon$  factor approx by running  $k = O(1/\epsilon^2)$  FM algo's iid.

Algo to get  $1 \pm \epsilon$  approx.

Bottom-k Algo:

- keep  $z_1, \dots, z_k := 1$  (counters.)

- @ item  $i$  in stream, compute  $h(i)$ .

maintain  $z_1 < z_2 < \dots < z_k$  smallest  $k$  hash values seen in the stream.

- Estimator:  $\frac{k}{z_k}$ .

$k=2$



Goal: to prove  $\hat{d} = \frac{k}{z_k}$  is within factor  $1 \pm \epsilon$  of  $d$  with  $\geq 90\%$  prob.

Claim:  $\Pr_{\pi} [\hat{d} > d \cdot (1 + \epsilon)] \leq 0.05$  (for  $k = \Theta(1/\epsilon^2)$ ).

Pf:  $\hat{d} > d \cdot (1 + \epsilon) \Leftrightarrow \frac{k}{z_k} > d(1 + \epsilon)$

$\Leftrightarrow z_k < \frac{k}{d(1 + \epsilon)}$

$\Leftrightarrow$  there are at least  $k$  items (in stream) whose  $h(\cdot)$  is  $< \frac{k}{d(1 + \epsilon)}$

$\pi = \Pr$  [ choose  $d$  r.v.  $\in [0, 1]$ ,  
at least  $k$  of them  $< \frac{k}{d(1 + \epsilon)}$  ]

$X_i \equiv$  indic. r.v. of event  $i$ th item has  $h(\cdot) < \frac{k}{d(1 + \epsilon)}$ .

$\Rightarrow$  D. r. v.  $> 1$  ?



$$11 = \{ \dots, L < X_i \leq L+1, \dots \} \leftarrow$$

$$E[X_i] = P_i \left[ \begin{matrix} \nearrow \\ \searrow \end{matrix} \right] = \frac{k}{d(1+\epsilon)} \quad (\text{assuming } \epsilon < 1)$$

$$= P_i[X_i=1] \cdot 1 + P_i[X_i=0] \cdot 0$$

$$\text{Var}[X_i] \leq E[X_i^2] \leq \frac{k}{d(1+\epsilon)} < \frac{k}{d}$$

$$E[\sum X_i] = d \cdot \frac{k}{d(1+\epsilon)} = \frac{k}{1+\epsilon}$$

$$\text{Var}[\sum X_i] = d \cdot \text{Var}[X_i] = k$$

By Chebyshev:  $P_r \left[ \sum X_i - \frac{k}{1+\epsilon} > \sqrt{20k} \right] \leq \frac{1}{20} = 0.05$

$\sum X_i > \frac{k}{1+\epsilon} + \sqrt{20k}$  Set  $k$  s.t.  $\frac{k}{1+\epsilon} + \sqrt{20k} \leq k$

$\ll k$  for  $k = \frac{1}{\epsilon^2}$  for  $c$  large enough.

$\mathbb{P} = P_r[\sum X_i > k] \leq 0.05.$

$$k > \frac{k}{1+\epsilon} + \sqrt{20k}$$

$$k \Rightarrow k(1-\epsilon + \epsilon^2) + \sqrt{20k}$$

$$\left( \begin{array}{l}
 k(\varepsilon - \varepsilon') \Rightarrow \sqrt{20}k \\
 \sqrt{k} \Rightarrow \sqrt{20} / \varepsilon - \varepsilon^2 \Rightarrow k \Rightarrow \frac{\sqrt{20}}{\varepsilon^2} \\
 \text{using } \varepsilon < 1/2.
 \end{array} \right.$$