

AA Lecture 6, 1/28/2021

HW1 due

HW2 on Course Works.

Heavy hitters:

f_x = frequency of x , $x \in [n]$.
 x is ϕ -HH if $f_x \geq \phi \cdot \sum_y f_y = \phi m$.

$$\text{Eg: } n = 2^{32}$$

CountMin Algo: linear sketch of f_x : $A \cdot f$.

of size: $O(\frac{\lg n}{\epsilon \phi})$ words

→ returns set $H \subseteq [n]$ s.t., with Prob. $\geq 1 - \frac{1}{k}$

• if x is ϕ -HH then $x \in H$

• if x is not $(1-\epsilon)\phi$ -HH, $x \notin H$.

Morris: can do approx $O(\lg \lg n)$ bits

• estimator: estimate \hat{f}_x for any given x .

to output H_s for each $x \in [n]$

$$\hat{f}_x = \min_i S_i[h_i(x)]$$

if $\hat{f}_x > \dots$

add x to H .

Output H .

Time: at least n (universe of items)
 way too large. (Note: $n \gg m$)

Obs: $|H| \leq \frac{1}{(1-\epsilon)\phi}$

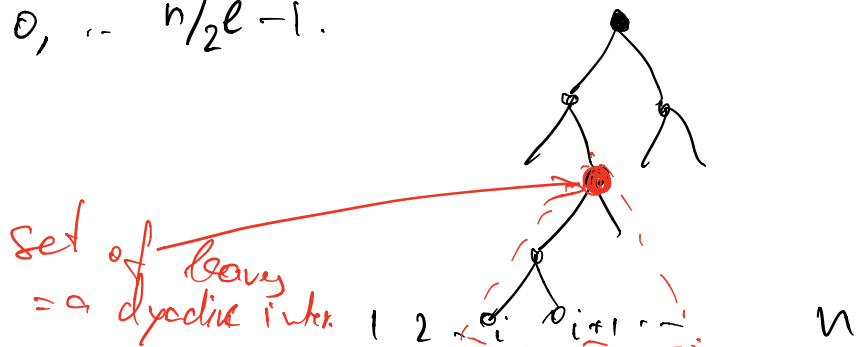
Goal: to compute ϕ -HH in time $\ll n$.

Thm: can get the same guarantees as Count Min (on H), using:

- $O\left(\frac{\lg^2 n}{\phi}\right)$ time to find HHS.
- $O\left(\frac{\lg^2 n}{\epsilon\phi}\right)$ space (in words).

Proof: n is a power of 2.

Def dyadic interval on $[n]$:
 an interval $[i \cdot 2^l, (i+1) \cdot 2^l]$
 $l = 0, 1, 2, \dots, \lg_2 n$
 $i = 0, \dots, n/2^l - 1$.



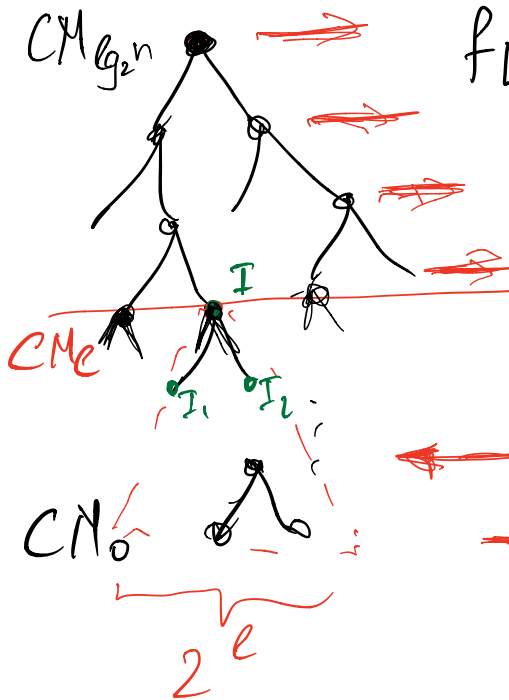
Fast Count Min:

store $\lg_2 n + 1$ sketches CM

CM_l : universe dyadic intervals @ level l

$[1, 2^l], [2^l+1, 2 \cdot 2^l], \dots$

$f_{[i \cdot 2^{l+1}, (i+1)2^{l+1}]} = \text{sum of frequencies of leaves.}$



each node has interval of $len=2^l$

items $[1, 2], [3, 4], \dots$
 $f_{[1,2]} = f_1 + f_2$

Like having $\lg_2 n + 1$ parallel streams.
 each with its own CM sketch
 (indep. of each other)

Obs: fix d. interval

$$I = I_1 \cup I_2$$

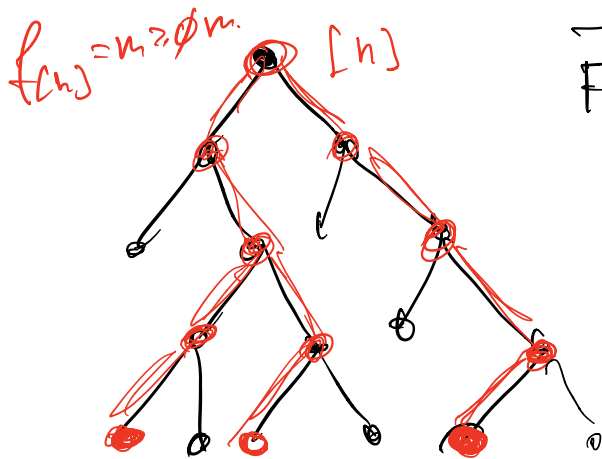
\uparrow level l \nwarrow \uparrow level $l-1$
 \dots

if node is HH

$f(I) = f(I_1) + f(I_2)$
 if I_1 or I_2 are ϕ -HHs \Rightarrow parent is HH.
 $\Rightarrow I$ is a ϕ -HH.

f_I for every level l $\sum f_I = m$.
 I \rightarrow dyad. interv. @ level l

if I_1 is ϕ -HH $\Rightarrow f_{I_1} \geq \phi \cdot m$
 $\Rightarrow f_I = f_{I_1} + f_{I_2} \geq \phi m$. \square



- Facts:
- 1) for every node, can check if HH
 - 2) in every level $\leq \frac{1}{(1-\epsilon)\phi}$ HHs.
 - 3) if node HH \Rightarrow parent is too.

Algo to find all HH:

HH(node v):

if v is leaf, add to H , return.

Let e_1 & e_2 be children.

\rightarrow use CM of e_1 .

if c_1 is HH, HH (c_1)

if c_2 is HH, HH (c_2).

Space: $O(\lg n)$ CM sketches $\Rightarrow O\left(\frac{\lg^2 n}{\epsilon \phi}\right)$.

Est time: $O(\lg n) \cdot 2 \cdot \frac{1}{(1-\epsilon)\phi}$ calls to ed .
 $\Rightarrow O\left(\frac{\lg^2 n}{(1-\epsilon)\phi}\right) = O\left(\frac{\lg^2 n}{\phi}\right)$.

Max frequency: $\|f\|_\infty = \max_x f_x$.

Other norms of f ?

$$\|f\|_1 = \sum_x f_x = m.$$

$$\|f\|_2 = \left(\sum_x f_x^2\right)^{1/2} \quad \leftarrow \text{TODAY.}$$

Frequency moments:

$$F_p = \sum_x f_x^p \quad F_2 = \|f\|_2^2.$$

Today: $F_2 / \|f\|_2$.

Example: $f^{(1)} = (1, 1, \dots, 1)$

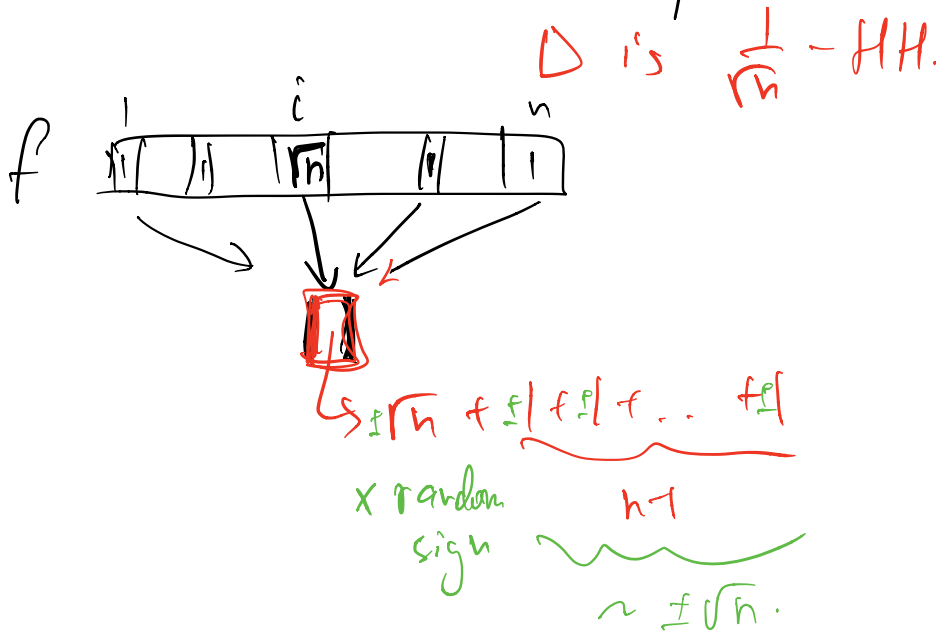
$$F_2 = n$$

vs $f^{(2)} = (1, 1, \dots, \sqrt{n}, \dots, 1)$.

$$F_2 = 2n - 1.$$

$\nearrow \sqrt{x^2}$

Can't distin. f^1 vs f^2 using CM. unless
space $\Theta(\sqrt{n})$.



Tug-of-War:

- init: $z = 0$

$\sigma: [n] \rightarrow \{\pm 1\}$ random.

- sketch: $z = \sum_x \sigma_x f_x$

@ seeing item x : $z := z + \sigma_x$

- estimator (for F_2): z^2

Analysis:

Claim 1: $E[z^2] = F_2$.

$$\begin{aligned}
 \text{pf: } \mathbb{E}_\sigma [z^2] &= \mathbb{E} \left[\left(\sum_x \sigma_x f_x \right)^2 \right] \\
 &= \mathbb{E} \left[\sum_{x,y} \sigma_x \sigma_y f_x f_y \right] \\
 &= \sum_{x,y} \underbrace{\mathbb{E} [\sigma_x \sigma_y]}_{\substack{\text{if } x=y, \Rightarrow 1 \\ x \neq y, \Rightarrow 0}} \cdot f_x f_y \\
 &= \sum_x f_x \cdot f_x = F_2. \quad \square
 \end{aligned}$$

Claim 2: $\text{var} [z^2] \leq O(F_2^2)$.

$$\begin{aligned}
 \text{pf: } \text{Var} [z^2] &= \mathbb{E} [z^4] - \underbrace{(\mathbb{E} [z^2])^2}_{} \\
 &= \mathbb{E} \left[\left(\sum_x \sigma_x f_x \right)^4 \right] \\
 &= \sum_{x,y,p,q} f_x f_y f_p f_q \cdot \mathbb{E} [\sigma_x \sigma_y \sigma_p \sigma_q] \\
 &= \sum_x f_x^4 + 3 \sum_{x,y} f_x^2 f_y^2.
 \end{aligned}$$

$$\leq \left(\sum_x f_x^2 \right)^2 + 3 \left(\sum_x f_x^2 \right)^2$$

$$= 4 F_2^2. \quad \square$$

$$\mathbb{E}[z^2] = F_2$$

$$\text{Var}[z^2] \leq 4 F_2^2$$

By Chebyshev: $\Pr[|z^2 - F_2| > 3 \cdot \sqrt{4 F_2^2}] \leq 1/9.$

$$z^2 = F_2 \pm 6 F_2 \text{ with prob } \geq 8/9.$$

To get better concentration (for $\epsilon \epsilon$ factor approx)

just "repeat" $k = \Theta(1/\epsilon^2)$ times.

To W+! - keep $k = \Theta(1/\epsilon^2)$ counters

$$z_1, \dots, z_k$$

$$\text{- each } z_i = \sum_x \sigma_x^i f_x,$$

where $\sigma^i: [n] \rightarrow \{\pm 1\}$ iid random.

$$\text{- Estimator: } z^2 = \frac{1}{k} \sum_{i=1}^k z_i^2.$$

Claim 1: $\mathbb{E} [z^2] = F_2$.

Claim 2: $\text{Var} [z^2] \leq \frac{1}{k} \cdot 4F_2^2$.

By Chebyshev:

$$\Pr \left[|z^2 - F_2| \geq 3 \cdot \sqrt{\frac{4F_2^2}{k}} \right] \leq \frac{1}{9}.$$

if this is false

$$\Rightarrow z^2 = F_2 \pm \frac{6}{\sqrt{k}} F_2 = F_2 \cdot \left(1 \pm \frac{6}{\sqrt{k}} \right).$$

set $k = 36/\epsilon^2$.

$$\Rightarrow z^2 = (1 \pm \epsilon) F_2 \text{ with probability } \geq 8/9.$$