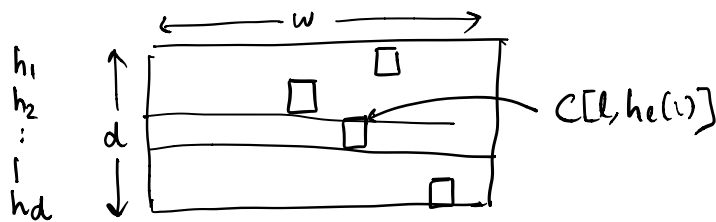Heavy Hitters
 — Count - Min
 — Count - Sketch

[Misra-Greis '82]    $f_i - \frac{m}{k} \leq \tilde{f_i} \leq f_i$    $m = F_1$

additive $\varepsilon F_1$ approximation for $k = \frac{1}{\varepsilon}$

Storage $\frac{1}{\varepsilon}$

Count-Min   [Cormode-Muthukrishnan]   CM-sketch

Array of counters width $w$ & depth $d$

For each row $i \in [d]$ : $h_i : [n] \longrightarrow [w]$



$C[\ell, h_\ell(i)]$

$h_1 \cdots h_d$ : pairwise independent hash fr $[n] \longrightarrow [w]$

For each element of stream
  $i \longleftarrow$ current element
  for $\ell = 1$ to $d$
    $C[\ell, h_\ell(i)] \longleftarrow C[\ell, h_\ell(i)] + 1$

Query $(i)$

$$\tilde{f_i} = \min_{\ell \in [d]} C[\ell, h_\ell(i)]$$

$$C[\ell, j] = \sum_{i : h_\ell(i) = j} f_i$$

<u>Analysis :</u>

$$\mathbb{E}[C[\ell, h_\ell(i)]] \leq f_i + \frac{m}{w}    m = F_1$$

Say $h_\ell(i) = b$

$$C[\ell, h_\ell(i)] = \sum_{i': h_\ell(i') = b} f_{i'}$$

$$\mathbb{E}[C[\ell, h_\ell(i)]] = f_i + \sum_{i' \neq i} \underbrace{\Pr[h_\ell(i') = b]} \cdot f_{i'}$$

$$= f_i + \frac{1}{w} \sum_{i' \neq i} f_{i'} \leq f_i + \frac{m}{w}$$

Also $C[\ell, h_\ell(i)] \geq f_i$

$$P\left[C[\ell, h_\ell(i)] \geq f_i + \frac{2m}{w}\right] = P\left[\colorbox{yellow}{$C[\ell, h_\ell(i)] - f_i \geq \frac{2m}{w}$}\right]$$

$$\leq \frac{\mathbb{E}[C[\ell, h_\ell(i)]] - f_i}{\frac{2m}{w}} \leq \frac{1}{2}$$

$h_1 \cdots h_d$ are independent

$$P\left[\tilde{f_i} \geq f_i + \frac{2m}{w}\right] = \prod_{\ell \in [d]} P\left[C[\ell, h_\ell(i)] \geq f_i + \frac{2m}{w}\right]$$

$$\leq \left(\frac{1}{2}\right)^d$$

$$w = \frac{2}{\varepsilon} \Rightarrow \frac{2m}{w} = \varepsilon m$$

$$d = \log_2 \frac{1}{\delta}, \quad \left(\frac{1}{2}\right)^d = \delta$$

$$\Pr\left[\tilde{f_i} \geq f_i + \varepsilon m\right] \leq \delta$$

$$f_i \leq \tilde{f_i} \leq f_i + \varepsilon m \qquad \text{w. prob. } 1-\delta$$

Space usage $\quad O\left(\frac{1}{\varepsilon} \log_2 \frac{1}{\delta}\right)$

Count - Sketch $\quad$ [C, Chen, Farach-Colton '02]

$h_1 \ldots h_d$    pairwise independent hash $f^n$ $[n] \rightarrow [w]$

$S_1 \ldots S_d$    ————"———— $[n] \rightarrow \{\pm 1\}$

For each element of stream

    $i \longleftarrow$ current elt.

    for $\ell = 1$ to $d$

$$C[\ell, h_\ell(i)] \longleftarrow C[\ell, h_\ell(i)] + S_\ell(i)$$

Query $(i)$

$$\hat{f}_i = \underset{\ell = 1 \text{ to } d}{\text{median}} \left\{ C[\ell, h_\ell(i)] \cdot S_\ell(i) \right\}$$

Analysis:

    Fix $i \in [n]$

$$Z_\ell = C[\ell, h_\ell(i)] \cdot S_\ell(i)$$

For $i' \in [n]$   $Y_{i'}$ indicator $\begin{cases} 1 & \text{if } h_\ell(i') = h_\ell(i) \\ 0 & \text{otherwise} \end{cases}$

$$\mathbb{E}[Y_{i'}^2] = \mathbb{E}[Y_{i'}] = \frac{1}{w}$$

$$Z_\ell = C[\ell, h_\ell(i)] \cdot S_\ell(i) = S_\ell(i) \sum_{i'} Y_{i'} \cdot f_{i'} \cdot S_\ell(i')$$

$$\mathbb{E}[Z_\ell] = f_i + \sum_{i' \neq i} \underbrace{\mathbb{E}[S_\ell(i) S_\ell(i') \cdot Y_{i'}]}_{} f_{i'}$$

                  $\mathbb{E}[S_\ell(i') S_\ell(i)] = 0$   for $i' \neq i$

                  $Y_{i'}$ independent of $S_\ell(i), S_\ell(i')$

$$\mathbb{E}[Z_\ell] = f_i$$

$$\text{Var}[Z_\ell] = \mathbb{E}\left[ \left( \sum_{i' \neq i} S_\ell(i) S_\ell(i') \cdot Y_{i'} f_{i'} \right)^2 \right]$$

$$= \mathbb{E}\left[ \sum_{i' \neq i} f_{i'}^2 Y_{i'}^2 + \sum_{\substack{i' \neq i'' \\ i', i'' \neq i}} f_{i'} f_{i''} S_\ell(i') S_\ell(i'') Y_{i'} Y_{i''} \right]$$

$$= \sum_{i' f_i} f_i^2 \; \mathbb{E}[y_i^2]$$

$$\leq \frac{\|f\|_2^2}{w}$$

$$w = \frac{3}{\varepsilon^2}$$

$$\Pr\left[|z_e - f_i| \geq \varepsilon \|f\|_2\right] \leq \frac{\text{Var}[z_e]}{\varepsilon^2 \|f\|_2^2} \leq \frac{1}{\varepsilon^2 w} \leq \frac{1}{3}$$

Now via Chernoff bound

$$\Pr\left[|\text{median}\{z_1 \ldots z_d\} - f_i| \geq \varepsilon \|f\|_2\right] \leq e^{-cd} \leq \delta$$

by choosing $d = O(\log(\frac{1}{\delta}))$

Space: $O(\frac{1}{\varepsilon^2} \log(\frac{1}{\delta}))$

| Comparison | Guarantee | Space |
|---|---|---|
| Misra-Greis | $f_i - \varepsilon\|f\|_1 \leq \tilde{f}_i \leq f_i$ | $\frac{1}{\varepsilon}$ |
| Count-Min | $f_i \leq \tilde{f}_i \leq f_i + \varepsilon\|f\|_1$ <br> w. prob. $1-\delta$ | $O(\frac{1}{\varepsilon} \log(\frac{1}{\delta}))$ |
| Count-Sketch | $|f_i - \tilde{f}_i| \leq \varepsilon\|f\|_2$ <br> w. prob. $1-\delta$ | $O(\frac{1}{\varepsilon^2} \log(\frac{1}{\delta}))$ |

$$\|f\|_1 \geq \|f\|_2$$

all 1's     $n$     $\sqrt{n}$

Heavy tailed   $f_i \sim \frac{1}{\sqrt{i}}$

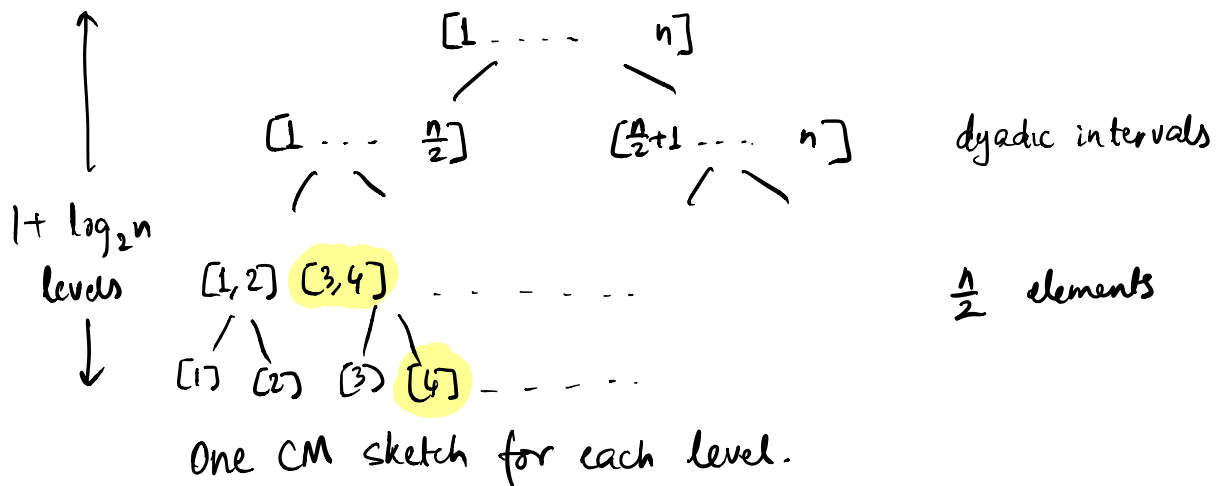$$\|f\|_1 = \theta(\sqrt{n}) \qquad \|f\|_2 = \theta(1)$$

We solved: point queries
estimate $f_i$ to within $\pm \varepsilon \|f\|_1$
$$\pm \varepsilon \|f\|_2$$

Heavy Hitters:

Output $L$ such that

$f_i \geqslant \varepsilon \|f\|_1 \Rightarrow i \in L$

$f_i < \frac{\varepsilon}{2} \|f\|_1 \Rightarrow i \notin L$

$\frac{\varepsilon}{2}\|f\|_1 \qquad \varepsilon\|f\|_1$



dyadic intervals

$[1 \cdots n]$

$[1 \cdots \frac{n}{2}] \qquad [\frac{n}{2}+1 \cdots n]$

$1 + \log_2 n$ levels

$[1,2] \quad [3,4] \quad \cdots$

$[1] \quad [2] \quad [3] \quad [4] \quad \cdots$

$\frac{n}{2}$ elements

One CM sketch for each level.

$\alpha$–heavy hitters w. failure prob. $\leq \delta$

each CM sketch has error parameter $\varepsilon = \frac{\alpha}{4}$

failure prob. $\eta = \frac{\delta}{(\log n)} \frac{\alpha}{4}$

At each level $j$ of tree, track $L_j$: heavy hitters at level $j$

$L_j$: contain all $\alpha$–heavy hitters at its level
no one below $\frac{\alpha}{2}$ heavy

For each of 2 children, point query child using CM sketch at level $j+1$

If child has $point\ query \geq \left(\frac{3\alpha}{4}\right) \|x\|_1$

include in $L_{j+1}$

$L$: list at leaf level.

Conditioned on correctness (no failure)

$L_j$ has size $\leq \frac{2}{\alpha}$

we query $\leq \frac{4}{\alpha}$ intervals at level $j+1$
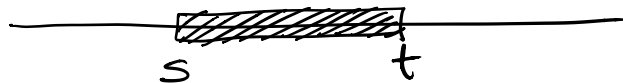
At most $Q \leq \frac{4}{\alpha} \log n$

We chose parameters so that each CM sketch has

failure prob $\leq \frac{\delta}{Q} = \frac{\delta \cdot \alpha}{4 \log n}$

union bound, failure prob. $\leq \delta$

Interval queries.

$\sum_{i \in [s,t]} f_i$



interval queries via point queries:

query time linear in interval size

error scales linearly in interval size

Instead:

Use interval based data structure

Each interval $[s,t]$ is union of at most $2 \log n$ dyadic intervals

Now query time $= O(\log n)$ (Query time of CM sketch)

error $\leq O(\log n) \cdot$ (error of CM-sketch at each level)