

## Heavy Hitters

stream of elts  $a_1, a_2, \dots, a_t, \dots$  e.g. IP addresses  
each  $a_i \in [n]$  ex:  $n = 2^{64}$

For any time  $t \leq n$

let  $f_x^t$ : # times elt  $x$  appears in  $a_1, a_2, \dots, a_t$  (frequency)

Goal: return all elts  $x$  s.t.  $f_x$  large  $> \frac{n}{k}$

using  $o(n, t)$  space

If  $f_x^t > \frac{n}{k}$ , add  $x$  to list of "heavy hitters"

Provably impossible with sublinear space

Modified goal:

① If  $f_x > \frac{n}{k}$  add it to list.

② If  $x$  added to list, then w.p.  $\geq 1 - \delta$ ,  $f_x^t \geq \frac{n}{k} - \epsilon n$

Ex Suppose  $k=20$ ,  $\epsilon=0.01$ ,  $\delta=\frac{1}{2^{10}}$

① If  $f_x > \frac{n}{20}$ ,  $x$  is definitely on list of Ht  
 $= 0.05n$

② if  $x$  added to list, then with prob  $\geq 1 - \frac{1}{2^{10}}$

$f_x > \frac{n}{k} - \epsilon n = 0.05n - 0.01n = 0.04n$

Modified goal:

① If  $f_x > \frac{n}{k}$  add it to list.

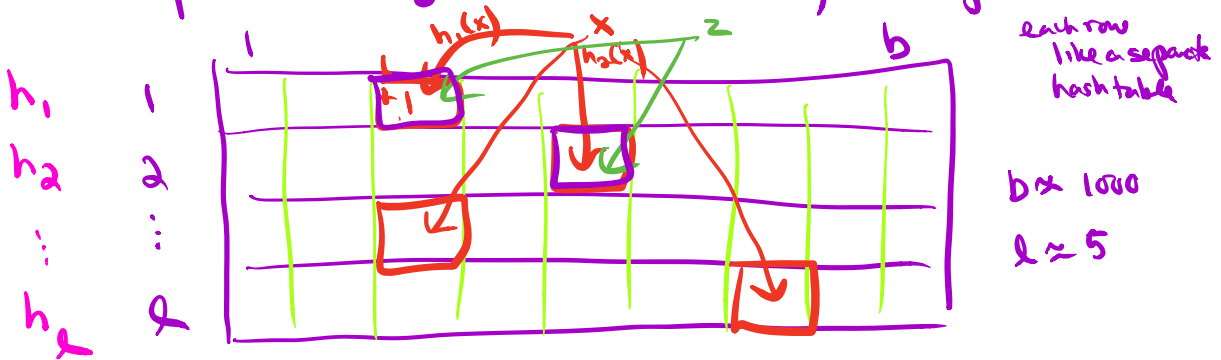
② If  $x$  added to list, then w.p.  $\geq 1 - \delta$ ,  $f_x^t \geq \frac{n}{k} - \epsilon n$

use elegant, very simple data structure Count-Min-Sketch.

# Count Min Sketch

designer specifies  $k, d, \epsilon$   
 $\Rightarrow b, \ell$ .

keep 2D array:  $\ell$  hash tables, each of size  $b$



when elt  $x$  shows up

$\text{Inc}(x)$ :  $\forall 1 \leq j \leq \ell$  increment  $\text{CMS}[j][h_j(x)]$

Observation 1: at any time  $t$   
 $\forall j, \forall x$   $\text{CMS}[j][h_j(x)] \geq f_x^t$

$\text{Count}(x)$ : return  $\min_{1 \leq j \leq \ell} \text{CMS}[j][h_j(x)]$

by Observation 1  $\text{Count}(x) \geq f_x^t$  if  $\text{Count}(x) > \frac{n}{k}$   
 add  $x$  to HT list.

## Assumptions:

① hash fns behave like random maps in the following sense  
 $\forall x, y \in [n]$   $\forall 1 \leq j \leq \ell$   $\Pr(h_j(x) = h_j(y)) = \frac{1}{b}$   
 $x \neq y$

② hash fns  $1 \leq j \leq \ell$  are indep of each other

## Analysis

Fix time  $t$ , elt  $x$

Define

$$Z_j = \text{CMS}[j][h_j(x)]$$

$$Z_j \geq f_x^t$$

$$Z_j = f_x^t + \sum_{\substack{y \neq x \\ y \in [n]}} f_y^t W_{xy}^+$$

$$W_{xy}^+ = \begin{cases} 1 & \text{if } h_j(x) = h_j(y) \\ 0 & \text{otherwise} \end{cases}$$

$$E(Z_j) = f_x^+ + \sum_{y \neq x} f_y^+ E(W_{xy}^+) \leq f_x^+ + \frac{1}{b} \leq f_x^+ + \frac{n}{b}$$

(0 a.w.)

$\stackrel{Pr(h_j(x)=h_j(y))}{=} \frac{1}{b}$

$$E(Z_j - f_x^+) \leq \frac{n}{b}$$

By Markov's Ineq ( $\alpha=2$ )

$$Pr(Z_j - f_x^+ \geq \frac{2n}{b}) \leq \frac{1}{2} \quad (*)$$

Markov's Inequality

X is a nonnegative r.v.

$$Pr(X > \alpha E(X)) \leq \frac{1}{\alpha}$$

$\uparrow$   
 $\alpha \geq 0$

$$Pr\left(\min_{1 \leq j \leq k} Z_j \geq f_x^+ + \frac{2n}{b}\right) \leq \frac{1}{2^k}$$

$$= Pr\left(Z_1 \geq f_x^+ + \frac{2n}{b}, Z_2 \geq f_x^+ + \frac{2n}{b}, \dots, Z_k \geq f_x^+ + \frac{2n}{b}\right)$$

$$= Pr\left(Z_1 \geq f_x^+ + \frac{2n}{b}\right) Pr\left(Z_2 \geq f_x^+ + \frac{2n}{b}\right) \dots Pr\left(Z_k \geq f_x^+ + \frac{2n}{b}\right)$$

$\stackrel{\text{indep of hashing}}{=} \text{indep of } Z_j\text{'s}$

by (\*)  $\leq \frac{1}{2^k}$

$$Pr\left(\min_{1 \leq j \leq k} Z_j \geq f_x^+ + \frac{2n}{b}\right) \leq \frac{1}{2^k}$$

$k, \epsilon, \delta$

Modified goal:

① If  $f_x > \frac{n}{k}$  add it to list.

② If x added to list, then w.p.  $\geq 1-\delta$ ,  $f_x^+ \geq \frac{n}{k} - \epsilon n$

$$\delta = \frac{1}{2^k} \quad \epsilon n = \frac{2n}{b}$$

$$k = \log\left(\frac{1}{\delta}\right) \quad \epsilon = \frac{2}{b}$$

$$b = 200$$

$$l = 5 \Rightarrow \epsilon = \frac{2}{b} = \frac{2}{200} = 0.01$$

$$l = 5 \Rightarrow \text{error prob is } \frac{1}{2^5} = \frac{1}{32}$$

$$\frac{n}{k} = 0.04n$$

✓ if  $x$  has  $f_x > 0.04n$   $x$  definitely output

if output  $x$ , then with prob  $\geq \frac{31}{32}$ ,  $f_x \geq \frac{n}{k} - \epsilon n$

$$= 0.04n - 0.01n$$

$$= 0.03n$$

$p$  prime number  $> n$

$$\mathcal{H} = \left\{ h(x) = (ex + g) \bmod p \bmod b \mid \begin{array}{l} 1 \leq e \leq p-1 \\ 0 \leq g \leq p-1 \end{array} \right\}$$

$$|\mathcal{H}| = p(p-1)$$

universal family.

$$h(x) = (e_1 x + g_1) \bmod p \bmod b \rightarrow e_1, g_1$$

$$\rightarrow e_2, g_2$$

$$\rightarrow e_3, g_3$$



If  $h$  is chosen uniformly at random from  $\mathcal{H}$

$$\forall x \neq y$$

$$\Pr(h(x) = h(y)) = \frac{1}{b}$$

$$\uparrow$$

$$h \in \mathcal{H}$$

## Markov's Inequality

$X$  nonnegative r.v. . For any positive constant  $c$ ,

$$\Pr(X \geq cE(X)) \leq \frac{1}{c}$$

Proof:

$$\begin{aligned} E(X) &= \sum_{x < cE(X)} x \Pr(X=x) + \sum_{x \geq cE(X)} x \Pr(X=x) \\ &\geq 0 + \sum_{x \geq cE(X)} cE(X) \Pr(X=x) \\ &= cE(X) \Pr(X \geq cE(X)) \end{aligned}$$

$$\Rightarrow \Pr(X \geq cE(X)) \leq \frac{1}{c}$$