

## MCMD L9 : Streaming-Counting Distinct Elements

### Streaming Algorithms

Stream :  $A = \langle a_1, a_2, \dots, a_m \rangle$   
 $a_i \in [n]$  size  $\log n$   
Compute  $f(A)$  in  $\text{poly}(\log m, \log n)$  space

-----

Flajolet + Martin '85  
Alon, Matias, Szegedy '99

$f_j = |\{a_i \in A \mid a_i = j\}|$

Goal:  $F_0 = |\{j \in [n] \mid f_j \geq 0\}|$   
number of distinct elements

$\text{zeros}(p) = \max\{i \mid 2^i \text{ divides } p\}$

#####

Init:

Choose random hash  $h : [n] \rightarrow [n]$

$z := 0$

Stream: A

if  $(\text{zeros}(h(a_i)) > z)$  then  $z := \text{zeros}(h(a_i))$

Output:  $2^{\{z+1/2\}}$

#####

----

How to implement  $h : [R] \rightarrow [m]$ ?

$h_a$  in H. Let  $a$  be a large real number: say  $a = \text{Unif}(0,1) * 10000$

Define  $h_a(x) = \text{floor}(m * \text{frac-part}(x*a))$

$\text{frac-part}(15.324) = 0.234$

$\text{floor}(15.324) = 15$

----

Let there be  $k$  distinct elements.

- we don't know answer, but used in analysis

Expect  $1/k$  distinct elements to have  $\text{zeros}(a_i) \geq \log k$

Expect no elements to have  $\text{zeros}(a_i) \gg \log k$

-----

Let  $X_{r,j}$  == indicator random variable for  $[\text{zeros}(h(j)) > r]$

$$Y_r = \sum_{\{j \text{ s.t. } a_i=j\}} X_{r,j}$$

Let  $t = z$  at end of stream.

$$Y_r > 0 \iff t \geq r$$

$$Y_r = 0 \iff t < r$$

$$E[X_{r,j}] = \Pr[\text{zeros}(h(j)) \geq r] = \Pr[2^r \text{ divides } h(j)] = 1/2^r$$

$$E[Y_r] = \sum_{\{j \text{ s.t. } a_i=j\}} E[X_{r,j}] = k/2^r$$

$$\begin{aligned} \text{Var}[Y_r] &= \sum_{\{j \text{ s.t. } a_i=j\}} \text{Var}[X_{r,j}] && (= E[(X_{r,j})^2] - E[X_{r,j}]^2) \\ &\leq \sum_{\{j \text{ s.t. } a_i=j\}} E[X_{r,j}^2] \\ &= \sum_{\{j \text{ s.t. } a_i=j\}} E[X_{r,j}] \\ &= k/2^r \end{aligned}$$

+++++

Markov Inequality

X a rv and  $a > 0$

$$\Pr[|X| \geq a] \leq E[|X|]/a$$

+++++

+++++

Chebyshev's Inequality:

Y a rv and  $b > 0$

$$\Pr[|Y - E[Y]| \geq b] \leq \text{Var}(Y)/b^2$$

+++++

using MI with  $X = (Y - E[Y])^2$  and  $a = b^2$

-----

$$\text{MI : } \Pr[Y_r > 0] = \Pr[Y_r \geq 1] \leq E[Y_r]/1 = k/2^r \quad (\text{E1})$$

given  $r < \log k$  then

$$\begin{aligned} \text{CI : } \Pr[Y_r = 0] &= \Pr[|Y_r - E[Y_r]| \geq k/2^r] \\ &\leq \text{Var}[Y_r]/(k/2^r)^2 \\ &\leq 2^r/k \end{aligned} \quad (\text{E2})$$

Algorithm output:  $\hat{k}$

$$\hat{k} = 2^{\lceil t+1/2 \rceil}$$

Let  $a$  == smallest integer s.t.  $2^{a+1/2} \geq 3k$ .

$$\Pr[\hat{k} > 3d] = \Pr[t \geq a] = \Pr[Y_a > 0] \leq k/2^a \leq \sqrt{2}/3 < 1/2$$

Let  $b =$  largest integer s.t.  $2^{\lfloor b+1/2 \rfloor} < k/3$ .

$\Pr[\hat{k} \leq d/3] = \Pr[t \leq b] = \Pr[Y_{b+1} = 0] \leq 2^{b+1}/k \leq \sqrt{2}/3 < 1/2$

( $\epsilon=3, \delta=1/2$ )-approximation

-----  
Median Trick

(make  $\delta$  arbitrary small)

Run  $s$  parallel, independent hash functions on the above procedure.

output:  $\hat{K} = \{ \hat{k}_1, \hat{k}_2, \dots, \hat{k}_s \}$

let  $\bar{k} = \text{median}[\hat{K}]$

$\bar{k} > 3k$  only if  $s/2$  values in  $\hat{K} > 3k$ .

Each  $\leq 3k$  w.p.  $1/2$  -- all independent

$$1/2^{\lfloor s/2 \rfloor} \leq \delta \quad (\text{where we choose } \delta)$$

solve for  $s$  :

$$2^{\lfloor s/2 \rfloor} \geq 1/\delta$$

$$s/2 \geq \log(1/\delta)$$

$$s \geq 2 \log(1/\delta)$$

Similar for lower bound:  $\delta \rightarrow \delta/2$

Using  $s = 2 \log(2/\delta)$ , take median  $\bar{k}$  is an ( $\epsilon=3, \delta$ )-approximation of # distinct elements.

$O(\log \log n)$  bits to store  $t$

$O(\log(1/\delta))$  hash functions

So:  $O(\log(1/\delta) * \log \log n)$  right?

oops, forgot to store hash function:

$O(\log n)$  bits to store hash function

So:  $O(\log(1/\delta) * \log n)$

-----

Better algorithm:

Space:  $O(\log m + (1/\epsilon^2) (\log(1/\epsilon) + \log \log m))$

( $\epsilon, \delta$ )-approximation

Hashes to smaller number of bins

Takes average to drive  $\epsilon$  down