

Today

- Wrap up AMS sketch for  $F_2$
- Distinct elements

Setting

Stream of data from  $X$

$f(x) = \text{frequency of } x = \# \text{ times } x \text{ appears in the stream}$

$$m = \text{length of the stream} \quad | \quad F_2 = \sum f^2(x)$$

Last time

$h(\cdot) : X \rightarrow \{-1, 1\} \leftarrow \text{fully random}$

$$y = \sum_{x \in X} h(x) \cdot f(x)$$

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

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

[How can we use this?]

Chebyshov's Inequality

$X$  - random variable, finite expectation & variance

$$\Pr(|X - \mathbb{E}[X]| \geq a \sqrt{\text{Var}[X]}) \leq \frac{1}{a^2} \quad \text{for any } a > 0$$

$k$  independent copies:  $y_1, y_2, y_3, y_4, \dots, y_k$

$$\text{Output } Z = \frac{\sum_{i=1}^k y_i^2}{k}$$

$\mathbb{E}[Z] = F_2$  independent variables

$$\text{Var}[Z] = \sum_i \text{Var}\left[\frac{y_i^2}{k}\right] = \frac{1}{k^2} \sum_i \text{Var}[y_i^2] \leq \frac{k}{k^2} \cdot 2F_2^2$$

$$= \frac{2F_2^2}{k}$$

Set  $k = \lceil 18/\varepsilon^2 \rceil$

$$\text{Var}[Z] \leq \frac{\varepsilon^2 F_2^2}{9}$$

Chebyshov's inequality gives:

$$\begin{aligned} \Pr(|Z - F_2| \geq \varepsilon F_2) &\leq \Pr(|Z - \mathbb{E}[Z]| \geq 3\sqrt{\frac{\varepsilon^2 F_2^2}{9}}) \\ &\leq \Pr(|Z - \mathbb{E}[Z]| \geq 3\sqrt{\text{Var}[Z]}) \leq \frac{1}{3^2} \end{aligned}$$

In other words:

W.p. at least  $\frac{8}{9}$ :

$$(1-\varepsilon)F_2 \leq Z \leq (1+\varepsilon)F_2$$

This is called  $(1+\varepsilon)$ -multiplicative approximation

Implementation of each  $y_i$ :

- start from empty counter
- for  $x$  in the stream add  $h(x)$

Space usage:  $O(1/\varepsilon^2)$  counters

Hash functions: 4-wise independence

Suffices to get  
all expectations right

such as  $E[h(x)h(y)h(z)h(t)]$

For  $X = [n]$ ,  $O(1)$  words of space each

How: ~~Discussion~~ Discussion tomorrow  
(bonus: secret sharing)

## Improve probability of success

$$\delta/g \rightarrow 1 - \delta$$

- Run  $O(\log(1/\delta))$  independent copies
- Return median of results

Why this works: Homework 1 (most likely)

## Distinct elements

Goal: Compute  $(1+\varepsilon)$ -multiplicative approximation to

$$F_0 = |\{x \in X : f(x) \neq 0\}|$$

## Algorithm

$h$  = random hash function from  $X$  to  ~~$\mathbb{N}$~~   $N$  s.t.

$$\Pr(h(x) = i) = 2^{-(i+1)}$$

$h(x) =$	0	1	2	3	$\dots$
$\Pr(\checkmark)$	$1/2$	$1/4$	$1/8$	$1/16$	$\dots$

Initially:

$$z \leftarrow 0$$

$$A \leftarrow \emptyset$$

For each element  $x$  in the stream:

if  $h(x) \geq z$

$$A \leftarrow A \cup \{(x, h(x))\}$$

while  $|A| \geq c \varepsilon^2$ :

$$z \leftarrow z + 1$$

remove from  $A$  all pairs  $(y, g)$   
s.t.  $g < z$

some large  
constant  ~~$c \geq 1000$~~   
 $(c > 576)$

$$\text{Output: } |A| \cdot 2^z$$

## Analysis:

Random variables:

$Z$  - final value ~~of~~ of  $z$

$$X_{i,x} \xrightarrow[\text{O.O.W.}]{\text{fith}} (x) \geq i \quad (\text{for } x \in X, i \in N)$$

$$Y_i = \sum_{x:f(x) \geq i} X_{i,x} \quad (\text{for } i \in N)$$

Output of the algorithm:  $Y_Z \cdot 2^Z$

$$\text{Incorrect output: } |Y_Z \cdot 2^Z - F_0| > \varepsilon F_0$$

$$\left| Y_Z - \frac{F_0}{2^Z} \right| > \frac{\varepsilon F_0}{2^Z}$$

If ~~the~~  $Z=0$ :  $|A| = Y_0 = F_0$

We output exact value

$$\textcircled{*} = \Pr(\text{incorrect output}) = \sum_{\substack{i=1 \\ \text{integer}}}^{\infty} \Pr\left(\left| Y_i - \frac{F_0}{2^i} \right| > \frac{\varepsilon F_0}{2^i} \wedge Z=i\right)$$

Select  $s$  s.t.

$$\frac{12}{\varepsilon^2} \leq \frac{F_0}{2^s} < \frac{24}{\varepsilon^2}$$

If  $s \leq 1$ ,  $F_0 \ll 9/\varepsilon^2$  and  $Z=0 \Rightarrow$  algorithm outputs exact value

4-6

$$\textcircled{A} \leq \sum_{i=1}^{s-1} \Pr\left(\left|Y_i - \frac{F_0}{2^i}\right| > \frac{\varepsilon F_0}{2^i}\right) + \sum_{i=s}^{\infty} \Pr(Z = i)$$

△                                   □

$$\square = \Pr(Z \geq s) = \Pr(Y_{s-1} \geq c/\varepsilon^2)$$

$$\stackrel{\text{Markov's}}{\leq} \frac{\mathbb{E}[Y_{s-1}]}{c/\varepsilon^2} = \frac{F_0/2^{s-1}}{c/s^2} = \frac{2\varepsilon^2}{c} \cdot \frac{F_0}{2^s}$$

$$\nearrow \leq \frac{2\varepsilon^2}{c} \cdot \frac{24}{\varepsilon^2} = \frac{48}{c} \quad \begin{matrix} \swarrow \\ \uparrow \\ (\text{large } c) \end{matrix} \quad \frac{1}{12}$$

Selection of  $s$

$$\Delta \leq ?$$

$$\mathbb{E}[Y_i] = \frac{F_0}{2^i}$$

$$\text{Var}[Y_i] = \sum_{x:f(x) > 0} \text{Var}[X_{i,x}] \leq \sum_{x:f(x) > 0} \mathbb{E}[X_{i,x}^2]$$

$$= \sum_{x:f(x) > 0} \mathbb{E}[X_{i,x}] = \mathbb{E}[Y_i] = \frac{F_0}{2^i}$$

$$\Delta = \sum_{i=1}^{s-1} \Pr\left(|Y_i - \mathbb{E}[Y_i]| > \frac{\varepsilon F_0}{2^i}\right) = \sum_{i=1}^{s-1} \Pr\left(|Y_i - \mathbb{E}[Y_i]| > \varepsilon \sqrt{\frac{F_0}{2^i}} \sqrt{\text{Var}[Y_i]}\right)$$

$$\leq \sum_{i=1}^{s-1} \frac{2^i}{\varepsilon^2 F_0} < \cancel{\frac{2^s}{\varepsilon^2 F_0}} \leq \frac{1}{\varepsilon^2} \cdot \frac{\varepsilon^2}{12} = \frac{1}{12}$$

$\uparrow$   
Chebyshov's  
inequality

$\downarrow$   
Selection  
of  $s$

$$\textcircled{*} \leq \Delta + \square \leq \frac{1}{12} + \frac{1}{12} = \frac{1}{6}$$

Conclusion: The algorithm outputs a  $(1+\varepsilon)$ -multiplicative approximation w. p. at least  $5/6$

Space usage:  $O(1/\varepsilon^2)$  elements of  $X$  +  $O(1/\varepsilon^2)$  clean integers

$\underbrace{\quad}_{\text{can be reduced by hashing into a small range (2-wise independence is sufficient)}}$        $\underbrace{\quad}_{\text{at most } O(\log \log m) \text{ bits with high probability for further space savings (m = stream length)}}$

Optimal:  $O(1/\varepsilon^2 + \log n)$  bits for  $X = [n]$

Hash functions  $h$ ?

- 2-wise independence suffices, needed for  $\text{Var}[Y_i]$
- non-uniform distribution on some range?  
see Homework 1