

# Streaming Algorithms

- Defn: In data stream model, there's an input stream  $\tau =: \langle a_1, \dots, a_m \rangle$  whose elements are tokens  $a_i$  from the universe  $[n] = \{1, \dots, n\}$ .  
The CPU has space D; which is a rand-access memory.

Goal: Solve qns. about  $\sigma$  with  $D \leq O(\log m + \log n)$ .  
L<sub>0</sub> m or n is **HUGE**.

# 1) # distinct elements in $\sigma$

- Let  $d := \#\{a_1, \dots, a_m\}$ . Estimate  $d$  by  $\hat{d}$  s.t.  
 $d/3 \leq \hat{d} \leq 3 \cdot d$ , whp & using space  
only  $s \leq O(\log n)$ .

[Trivial:  $s = d \in [n]$ ] (like constant-many tokens!)

- Alon, Matias & Szegedy (1999): Their algo. uses  
a p.i.-hash function  $h : [n] \rightarrow [n]$ , that  
requires only log-space. And the valuation map  
 $v_2(h) := \max\{i : 2^i \mid h\}$ .

- Idea: Keep track of  $v_2(h(a_i))$ , for  $a_i \in \sigma$ .  
Whp one of the  $d$  elements has an  
 $h$ -image divisible by  $2^{\lfloor \lg d \rfloor} \approx d$ .

- Algo:

- 0) Choose rnd  $h$  from p.i.-hash-family;
- 1)  $z \leftarrow 0$ ;
- 2) for (token  $j$  in  $\sigma$ )
  - 3) If  $v_2(h(j)) > z$  then  $z \leftarrow v_2(h(j))$ ;
- 4) OUTPUT  $2^{z+0.5}$ ;

## Analyse:

- for each token  $j \in [n]$  &  $r \geq 0$ , define

$$\underline{X_{r,j}} := \begin{cases} 1, & \text{if } v_r(h(j)) \geq r \\ 0, & \text{else} \end{cases}$$

$$\& \underline{Y_r} := \sum \{ X_{r,j} \mid j \text{ s.t. } j \in \tau \}.$$

- Let  $\underline{T} :=$  terminal value of  $\underline{z}$  (when algo. stops).

$$\triangleright Y_r > 0 \iff T \geq r.$$

$$\triangleright Y_r = 0 \iff T \leq r-1.$$

$$\triangleright E[X_{r,j}] = P(2^r \text{ divides } h(j)) = \frac{1}{2^r}.$$

$$\triangleright E[Y_r] = d/2^r.$$

$$\triangleright \text{Var}(Y_r) = \sum_{j \in \sigma} \text{Var}(X_{r,j}) \leq \sum_{j \in \sigma} E[X_{r,j}^2]$$

(pri. hash)

$$= \sum_{j \in \sigma} E[X_{r,j}] = E[Y_r] = d/2^r.$$

- Output is  $\hat{d} := 2^{T+0.5}$ .

- { Let  $a \in \mathbb{Z}$  be the smallest :  $2^{a+0.5} \geq 3d$ .  
Let  $b \in \mathbb{Z}$  " largest :  $2^{b+0.5} \leq d/3$ .

$$\begin{aligned} \triangleright P(\hat{d} \geq 3d) &= P(T \geq a) = P(Y_a > 0) = P(Y_a \geq 1) \\ &\leq E[Y_a]/1 = d/2^a \leq \frac{\sqrt{2}}{3} < 1. \end{aligned}$$

$$\begin{aligned} \triangleright P(\hat{d} \leq d/3) &= P(T \leq b) = P(Y_{b+1} = 0) \\ &\leq P(|Y_{b+1} - E[Y_{b+1}]| \geq d/2^{b+1}) \leq \frac{\text{var}(Y_{b+1})}{(d/2^{b+1})^2} \\ &= 2^{b+1}/d \leq \sqrt{2}/3. \end{aligned}$$

$$\triangleright P\left(\frac{d}{3} \leq \hat{d} \leq 3d\right) \geq 1 - \frac{2\sqrt{2}}{3} > 0.$$

Boosting: The prob. can be made  $1 - 2^{-R(k)}$  by running  $k$  independent (& parallel) copies of the algo. for the stream  $\sigma$ . And OUTPUT the Median.

Chernoff. bd. gives

Jhm 1:  $P(d/3 \leq \text{output} \leq 3d) \geq 1 - 2^{-R(k)}$  &  $d \leq R \cdot \lg n$ .

## 2) Detecting a Heavy-hitter

- Want to check whether  $\sigma$  has element of "unusually high" frequency. (Perhaps, it's an attack on the server!)

- Let  $f_j :=$  frequency of  $j$  in  $\sigma$ .

$$\triangleright \sum_{j \in [n]} f_j = m =: \underline{F_1} \quad (\underline{\text{first moment}})$$

- Idea:  
• Compare  $F_1$  with the second moment  $F_2 := \sum_j f_j^2$ .  
•  $F_2 \geq 0.9 \times F_1^2 \Rightarrow$  Perhaps a heavy-hitter!

▷ Compute  $F_1$  by a counter in space  $\ell \text{ fm}$ .

Qn: How to compute  $F_2$ ?

- [AMS99] algo. to compute  $F_2$  uses a hash.fn.

$h: [n] \rightarrow \{-1, 1\}$  from 4-independent hash family.

(Exercise: Construct  $h$  by using cubic polynomials modulo a prime.)

Alg: 0) Pick rnd 4-indep. hash  $h: [n] \rightarrow \{\pm 1\}$ ;

1)  $Z \leftarrow 0$ ;

2) for (token  $j$  in  $\sigma$ )

    3)  $Z \leftarrow Z + h(j)$ ;

4) Output  $Y := Z^2$ ; // Is  $Y$  related to  $F_2$ ?

Analyse:  $\triangleright Z = \sum_{j \in [n]} f_j \cdot h(j)$ .

$$\triangleright E[Z] = \sum_j f_j \cdot E[h(j)] = 0.$$

$$\triangleright E[Y] = E\left[\left(\sum_j f_j \cdot h(j)\right)^2\right] = E\left[\sum_j f_j^2 + \sum_{i \neq j} f_i \cdot f_j \cdot h(i) \cdot h(j)\right]$$

$$= F_2 + \sum_{i \neq j} f_i \cdot f_j \cdot E[h(i)] \cdot E[h(j)] = F_2.$$

$\nwarrow$  2-wise indef.  $h$

$\Rightarrow$  We're on the right track!

$$\triangleright \text{var}(Y) = E[Y^2] - E[Y]^2 = E[Z^4] - E[Y]^2$$

$$= E\left[\left(\sum_j f_j \cdot h(j)\right)^4\right] - F_2^2$$

$$= E\left[\sum_j f_j^4\right] + 3 \cdot E\left[\sum_{i \neq j} f_i^2 f_j^2\right] - F_2^2$$

R 4-wise-indep. of  $h$

$$= \sum_j f_j^4 + 3 \cdot \sum_{i \neq j} f_i^2 f_j^2 - \sum_j (f_j^2)^2 - \sum_{i \neq j} f_i^2 f_j^2$$

$$= 2 \cdot \sum_{i \neq j} f_i^2 f_j^2 \leq 2 \cdot F_2^2$$

$$\triangleright \text{Chebyshew} \Rightarrow P(|Y - F_2| > \alpha \cdot F_2) \leq \frac{2F_2^2}{(\alpha F_2)^2} = 2/\alpha^2.$$

$\Rightarrow Y \text{ approx. } F_2 \text{ well!}$

Boosting: Reduce the variance by keeping  
k indep. estimators  $\{Y_1, \dots, Y_k\}$  &  
OUTPUT  $\underline{y'} := \frac{Y_1 + \dots + Y_k}{k}$ .

▷  $E[y'] = E[Y] = F_2$  &  
▷  $\text{var}(y') = \text{var}(Y)/k$ .

$\Rightarrow$  gives a much better approx. of  $F_2$ ,  
in space  $O(k \cdot lgm + k \cdot lgn)$ .