

## Frequency moments:

$a_1, \dots, a_T$  stream of elems from universe  $U = \{0, 1, \dots, p\}$

$f_j^+$  # of times  $j \in U$  appears in  $\{a_1, \dots, a_T\}$

$$f_j \triangleq f_j^+$$

$$F_k = k^{\text{th}} \text{ frequency moment of stream} = \sum_{j \in U} f_j^k$$

Algorithm for approximating  $F_2$  [Alon, Mathias, Szegedy]

- devise randomized unbiased estimator of  $F_2$

that can be computed on the fly

- repeat for accuracy

$h: U \rightarrow \{\pm 1\}$  equally likely

$Z := 0$

for  $i = 1$  to  $T$

$Z := Z + h(a_i)$  [ $++$  if  $h(a_i) = 1$      $--$  if  $h(a_i) = -1$ ]

Return  $X = Z^2$

Lemma

$$E(X) = F_a$$

Proof

$$Z = \sum_{j \in U} h(j) f_j$$

$$\begin{aligned} E(Z^2) &= E\left[\left(\sum_{j \in U} h(j) f_j\right)^2\right] = E\left[\sum_{j \in U} \boxed{h(j)^2} f_j^2 + \sum_{j \in U, k \in U} h(j) h(k) f_j f_k\right] \\ &= F_a \qquad \qquad \qquad E(h(j) h(k)) = 0 \end{aligned}$$

Note: so far only used pairwise independence

To show that  $Z^2$  is close to  $F_2$ , need to compute variance  
(and then can use trick of repetition to reduce variance)

Lemma:

$$\text{Var}(X) \leq 3F_2^2$$

$\Rightarrow$  repeating  $t = \frac{2}{\epsilon^2 \delta}$  times and averaging

$$\Rightarrow \text{Var}(\bar{X}) \leq \frac{\text{Var}(X)}{t} \leq \frac{\epsilon^2 \delta \text{Var}(X)}{2} \leq \epsilon^2 \delta F_2^2$$

$$\Rightarrow \Pr(|\bar{X} - F_2| > \epsilon F_2) \leq \frac{\text{Var}(\bar{X})}{\epsilon^2 F_2^2} \leq \delta$$

Space:  $O\left(\log m + \frac{\log T}{\epsilon^2 \delta}\right)$

[Actually can be improved to

Guarantee

$$\tilde{F}_2 (= \bar{X}) \in [F_2 \pm \epsilon F_2] \text{ with prob } \geq 1 - \delta$$

$$O\left(\log m + \frac{\log \log T}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$$

Proof of lemma:

$$E(Z^4) = E\left(\sum_{j \in U} h(j) f_j\right)^4$$

$$= \sum_{j_1} \sum_{j_2} \sum_{j_3} \sum_{j_4} E[h(j_1)h(j_2)h(j_3)h(j_4)] f_{j_1} f_{j_2} f_{j_3} f_{j_4}$$

$E[h(j_1)h(j_2)h(j_3)h(j_4)] = 0$  if  $\exists$  elt that appears only once

suppose  $j_1$  appears once

$$E[h(j_1)h(j_2)h(j_3)h(j_4) | h(j_2)h(j_3)h(j_4) = v] = 1 \cdot v + (-1) \cdot v = 0$$

$\Rightarrow$  only relevant terms are those where all terms appear an even # of times

$$\rightarrow = \sum_{j \in U} E[h(j)^4] f_j^4 + \binom{4}{2} \sum_{j_1 \neq j_2} E[h(j_1)^2] E[h(j_2)^2] f_{j_1}^2 f_{j_2}^2$$

$$= \sum_j f_j^4 + 6 \sum_{j_1 \neq j_2} f_{j_1}^2 f_{j_2}^2 \quad \ll 3 F_2^2$$

$$\text{since } F_2^2 = \left(\sum_j f_j^2\right) \left(\sum_j f_j^2\right) = \sum_j f_j^4 + 2 \sum_{j_1 \neq j_2} f_{j_1}^2 f_{j_2}^2$$

4-wise independence sufficient!

Summing up:

sublinear space    heavy hitters,  $F_0$  (#distinct elems),  $F_2$

randomization

necessary

necessary

approx

necessary

necessary

necessary