

## Estimating $F_2$ in Turnstile Streams

Recall  $F_2(\sigma) = \sum_i f_i^2 = \|f\|_2^2$ . In general,  $F_k(\sigma) = \sum_i f_i^k$ .

Recall we saw "AMS Sampling", which could  $(\epsilon, \delta)$ -approximate  $F_k$  using space  $O\left(\frac{n^{1-\frac{1}{k}} \log(n)}{\epsilon}\right)$  in insertion-only streams.

Here, we will see an algorithm that can  $(\epsilon, \delta)$ -approximate  $F_2$  in space  $O\left(\frac{(\log m) \log n}{\epsilon^2}\right)$  in turnstile streams (due to A McGregor, called the Tug-of-War Sketch).

### Basic Estimator.

- Choose a random hash function  $h: [n] \rightarrow \{-1, 1\}$  from a 4-wise independent hash family.
- Initialize  $x \leftarrow 0$ .
- When processing update  $(a_j, \delta_j)$   

$$x \leftarrow x + \delta_j \cdot h(a_j)$$
- Output  $x^2$ .

Same as Count-sketch with a single counter

Note:  $h$  takes  $\Theta(\log n)$  bits to store, and  $x$  takes  $\Theta(\log(M))$  bits.  
 Let  $X$  denote the random variable given by the ~~output~~ value of  $x$  at the end of the stream. Let  $Y_i$  be

Let  $X = \sum_{i=1}^n f_i \cdot Y_i$ , and the returned estimate is  $X^2$ .

Claim:  $E[X] = F_2(\sigma)$ .

Proof:  $E[X] = E\left[\left(\sum_i f_i Y_i\right)^2\right] = E\left[\sum_i f_i^2 Y_i^2 + \sum_{i \neq j} \sum_i f_i f_j Y_i Y_j\right]$

Linearity of expectation  $\rightarrow = \sum_i E[f_i^2 Y_i^2] + \sum_i \sum_{j \neq i} E[f_i f_j Y_i Y_j]$  Parameter dependence of  $h$   
 $= \sum_i f_i^2 + \sum_i \sum_{j \neq i} f_i f_j E[Y_i Y_j]$   $E[Y_i Y_j] = F_2$

Claim:  $\text{Var}[X^2] \leq 2 F_2^2$ .

Proof:  $\text{Var}[X^2] = \mathbb{E}[X^4] - \mathbb{E}[X^2]^2 = \mathbb{E}[X^4] - F_2^2$ . (\*)

$$\mathbb{E}[X^4] = \mathbb{E}\left[\left(\sum_i f_i Y_i\right)^4\right] = \mathbb{E}\left[\sum_i \sum_j \sum_k \sum_l f_i f_j f_k f_l Y_i Y_j Y_k Y_l\right]$$

(linearity of expectation)

$$= \sum_i \sum_j \sum_k \sum_l f_i f_j f_k f_l \mathbb{E}[Y_i Y_j Y_k Y_l]. \quad (**)$$

Any term in (\*) with ~~more than~~ of the indices  $(i, j, k, l)$  appearing exactly

~~once in the 4-tuple~~ is 0.

$$\begin{aligned} \text{E.g., if } i \notin \{j, k, l\}, \text{ then } \mathbb{E}[Y_i Y_j Y_k Y_l] \\ \text{4-wise independent} &= \mathbb{E}[Y_i] \mathbb{E}[Y_j Y_k Y_l] \\ &= 0 \end{aligned}$$

So ~~the~~ only non-zero terms are of the form

- one index occurring 4 times or
- two indices occurring twice each

$$\text{So } (*) = \sum_i f_i^4 \mathbb{E}[Y_i^4] + 6 \sum_{i \neq j} f_i^2 f_j^2 \mathbb{E}[Y_i^2 Y_j^2]$$

(4) permutations of  $(i, i, j, j)$

purpose of the  
rest of the calculation is to  
relate the expression to  $F_2^2$ .

$$\begin{aligned} &= \sum_i f_i^4 + 6 \left[ \sum_{i \neq j} f_i^2 f_j^2 \right] = \cancel{\sum_i f_i^4} + 3(F_2^2 - F_4) \\ &\quad = \frac{1}{2} \left( \sum_{i=1}^n f_i^2 \right)^2 - \sum_{i=1}^n f_i^4 = 3F_2^2 - 2F_4 \\ &\quad \leq 3F_2^2 \end{aligned}$$

$$\text{So } \text{Var}[X^+] \stackrel{\text{defn}}{=} E[X^4] - F_2^2 \leq 3F_2^2 - F_2^2 = 2 \cdot F_2^2$$

Do the Median-of-Means trick to obtain the final estimator. That is, average  $O(\frac{1}{\epsilon^2})$  copies of the basic estimator to drop its variance below (say)  $\frac{\epsilon^2 F_2^2}{4}$ .

Then Chebychev's inequality says that

$$\Pr[\text{estimate is off from expectation by more than } \epsilon \cdot F_2] \leq \frac{1}{4}.$$

So a median of  $O(\log(\frac{1}{\delta}))$  copies of the above is an  $(\epsilon, \delta)$ -approximation algorithm.