## Today: AMS sketch for $F_2$

Review: Frequency moments $F_p$

- Input/data: multiset of items from $X$
- Frequencies $f: X \to \mathbb{N}$

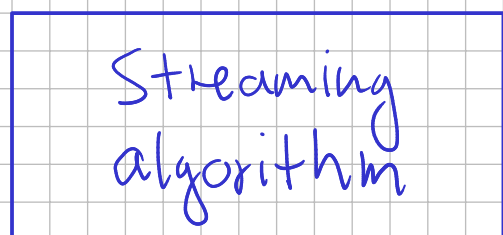$$f(x) = \#\text{occurrences of } x \in X$$

- p-th frequency moment

$p \in (0, \infty)$

$$F_p(\text{the data}) = \sum_{x \in X} |f(x)|^p$$

absolute value because in general $f(x)$ can be negative

---

## Streaming algorithms

big stream of data

┌─────────────┐
│ Streaming   │  ⟵  3, 11, 4, 3, 3, 1, 4, ...
│ algorithm   │
└─────────────┘

Today: we want to estimate $F_2(\text{stream})$ in small space

# Solution: AMS sketch
↑
Alon–Mattias–Szegedy (1996)

> This paper was very impactful.
> It inspired a lot of interest
> in the streaming model.

Basic estimator:

$$h : X \to \{-1, +1\} \text{ is "fully random"}$$
↑
each $h(x)$ independent
and uniformly distributed

Maintain $Y = \sum_{x \in X} h(x) f(x)$

Quick check: $\mathbb{E}[Y] = \sum_{x \in X} \underbrace{\mathbb{E}[h(x)]}_{=0} \cdot f(x) = 0$
↑
Hmm...

Not very
interesting
But...

$$\mathbb{E}[Y^2] = \mathbb{E}\left[\left(\sum_{x} h(x) f(x)\right)^2\right]$$

$$= \mathbb{E}\left[\sum_x \underbrace{h^2(x)}_{=1} f^2(x) + \sum_{\substack{x,y \\ x \neq y}} h(x)h(y)f(x)f(y)\right]$$

$$= \sum_x f^2(x) + \sum_{\substack{x,y \\ x \neq y}} f(x)f(y) \underbrace{\mathbb{E}\left[f(x)f(y)\right]}$$

$\nearrow$ *independence*

$$= \mathbb{E}\left[f(x)\right] \cdot \mathbb{E}\left[f(y)\right] = 0$$

$$= \sum_x f^2(x) = F_2 \leftarrow \text{exactly what we want}$$
$$\underline{\underline{\text{in expectation}}}$$

this scenario: $Y^2$ is an $\underline{\text{unbiased}}$ estimator for $F_2$

$$\left[\begin{array}{l} \text{Question: Is having an unbiased estimator} \\ \text{good enough?} \end{array}\right]$$

$$\text{Var}\left[Y^2\right] = \underbrace{\mathbb{E}\left[(Y^2)^2\right]}_{\boxed{\square}} - \underbrace{\left(\mathbb{E}\left[Y^2\right]\right)^2}_{\triangle}$$

$$\boxed{\square} = \mathbb{E}\left[(Y^2)^2\right] = \mathbb{E}\left[Y^4\right]$$

$$= \mathbb{E}\left[\sum_{x,y,z,t} h(x)h(y)h(z)h(t) f(x)f(y)f(z)f(t)\right]$$

$$= \sum_{x,y,z,t} \mathbb{E}\big[h(x)h(y)h(z)h(t)\big] \cdot f(x)f(y)f(z)f(t)$$

Up to four different elements $x, y, z, t$. If one of them occurs odd number of times, this expectation is 0.

Non-zero cases:

① $x = y = z = t$

② two pairs of different elements
(example: $x = t \neq y = z$)

$$\square = \sum_{x} f^4(x) + 3 \sum_{\substack{x,y \\ x \neq y}} f^2(x) f^2(y)$$

①          ②

$$\triangle = \Big(\mathbb{E}[Y^2]\Big)^2 = F_2^2 = \Big(\sum_{x} f^2(x)\Big)^2$$

$$= \sum_{x} f^4(x) + \sum_{\substack{x,y \\ x \neq y}} f^2(x) f^2(y)$$

$$\text{Var}[Y^2] = \boxed{\square} - \boxed{\triangle} = 2 \sum_{\substack{x, y \\ x \neq y}} f^2(x) f^2(y)$$

$$\leq 2 \sum_{x, y} f^2(x) f^2(y)$$

$$= 2 \left( \underbrace{\sum_x f^2(x)}_{F_2} \right) \left( \underbrace{\sum_y f^2(y)}_{F_2} \right) = 2 F_2^2$$

[How can we use this?]

Chebyshev's inequality

$X$ — random variable with finite expectation & variance

$$\Pr\left[ |X - \mathbb{E}[X]| \geq a \sqrt{\text{Var}[X]} \right] \leq \frac{1}{a^2}$$

for any $a > 0$

Use $k$ independent copies: $Y_1, Y_2, \ldots, Y_k$

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

$$\mathbb{E}[Z] = \frac{k \, \mathbb{E}[Y^2]}{k} = \mathbb{E}[Y^2] = F_2$$

$$\mathrm{Var}[Z] = \frac{1}{k^2} \mathrm{Var}\left[\sum_{i=1}^{k} Y_i^2\right] = \frac{1}{k^2} \cdot \sum_{i=1}^{k} \mathrm{Var}[Y_i^2]$$

<span style="color:red">↑ independent variables</span>

$$= \frac{k}{k^2} \mathrm{Var}[Y^2] \leq \frac{2 F_2^2}{k}$$

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

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

Via Chebyshev's inequality :

$$\Pr\left[|Z - F_2| \geq \varepsilon F_2\right] = \Pr\left[|Z - F_2| \geq 3\sqrt{\frac{\varepsilon^2 F_2^2}{9}}\right]$$

$$\leq \Pr\left[|Z - \mathbb{E}[Z]| \geq 3\sqrt{\mathrm{Var}[Z]}\right] \leq \frac{1}{3^2} = \frac{1}{9}$$

<span style="color:red">We will finish the discussion of AMS sketch next time</span>