

Today:

- A few more comments on the AMS sketch
- Amplifying the probability of a good estimate
- Chernoff bounds
- Constructing k -wise independent hash functions

} Next time

Review: AMS sketch for $F_2 = \sum_x f^2(x)$

Last time:

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

number of occurrences (aka. frequency) of x

where each $Y_i = \sum_{x \in X} h_i(x) f(x)$

hash function independently assigning each $x \in X$ to $\{-1, 1\}$ uniformly at random

Analysis last time:

If $k = \lceil 18/\epsilon^2 \rceil$, then

$$\Pr[|Z - F_2| \geq \epsilon F_2] \leq \frac{1}{9}$$

\Leftrightarrow equivalently

$$(1-\epsilon)F_2 \leq Z \leq (1+\epsilon)F_2 \quad \text{w.p. } \geq \frac{8}{9}$$

this kind of guarantee: multiplicative
 $(1+\epsilon)$ -approximation

Implementing in small space as a streaming algorithm:

Each Y_i :

- Initially, $Y_i = 0$

- For each insertion $x \in X$

$$Y_i \leftarrow Y_i + h_i(x)$$

For each deletion $x \in X$

$$Y_i \leftarrow Y_i - h_i(x)$$

$f(x)$ increases by 1

$f(x)$ decreases by 1

This keeps $Y_i = \sum_{x \in X} h_i(x) f(x)$

Outputting estimate:

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

Space usage (ignoring hash functions): $O(1/\epsilon^2)$

How about hash functions?

4-wise independence suffices

to make our analysis of

$\mathbb{E}[h(x)h(y)h(z)h(t)]$ valid

(and $\mathbb{E}[h(x)h(y)]$)

we get later to how to construct such functions in small space

How to make the probability of mult. $(1+\epsilon)$ -approximation higher?

Say, $\geq 1-\delta$ for $\delta \in (0, \frac{1}{9})$?

Compared to Count-Min Sketch (CMS):

- CMS: never underestimates, so take minimum of several basic estimators and it suffices that one of them is close

- AMS: can both overestimate and underestimate so taking maximum or minimum of several estimators won't work

Idea 1: Apply Chebyshev's inequality with higher number k of basic estimates we average

Last time via Chebyshev's: bound on variance of 2

$$\Pr \left[|Z - \underbrace{\mathbb{E}[Z]}_{= F_2}| \geq a \sqrt{\frac{2 F_2^2}{k}} \right] \leq \frac{1}{a^2}$$

for any $a > 0$

To bound the probability by δ ,

set $a = \sqrt{\frac{1}{\delta}}$:

$$\Pr \left[|Z - F_2| \geq \sqrt{\frac{2}{\delta k}} F_2 \right] \leq \delta$$

Next set $k = \lceil \frac{2}{\delta \epsilon^2} \rceil$, to get

$$\Pr \left[|Z - F_2| \geq \epsilon F_2 \right] \leq \delta$$

to get multiplicative $(1 + \epsilon)$ -approximation

Total space: $O\left(\frac{1}{s\epsilon^2}\right)$

Overhead:

- CMS: $\sim \log(1/s)$

- here: $\sim 1/s$

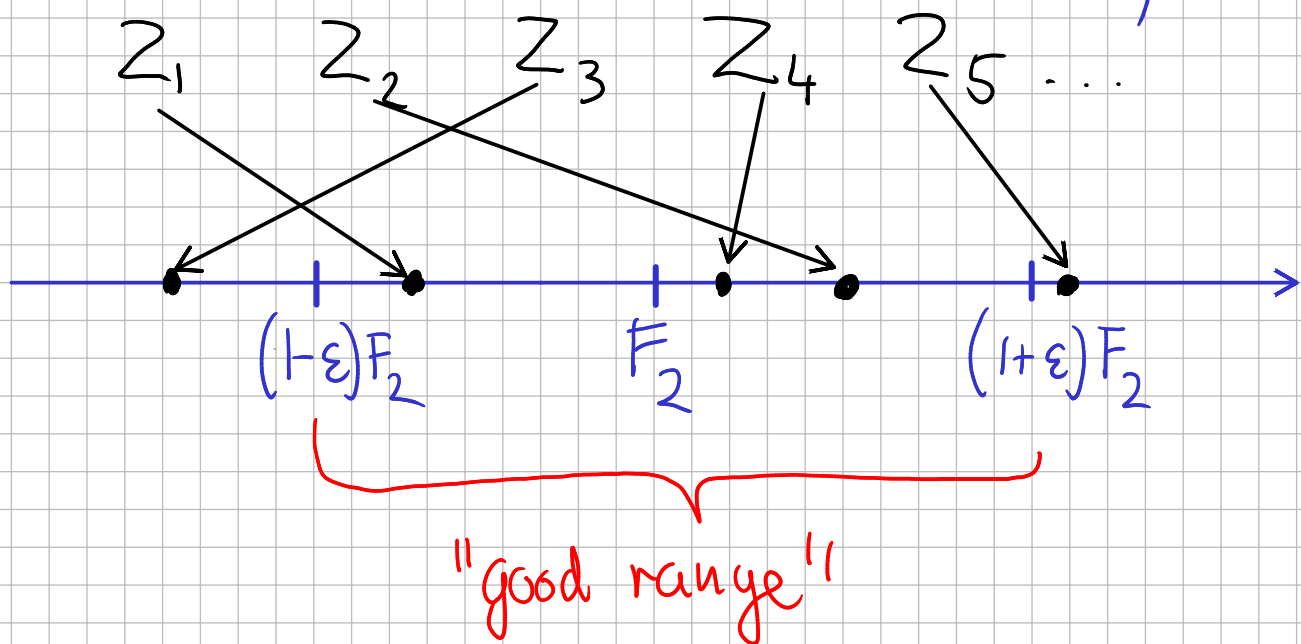
Can we do better?

Idea 2: Take the median of multiple estimators Z_i for $k = \lceil 18/\epsilon^2 \rceil$

Algorithm:

- Run k' independent Z_i 's (each outputs the mean of $k = \lceil 18/\epsilon^2 \rceil$ independent copies of the basic estimator Y^2)
to be set later
- Output the median of their estimates

Big picture: use "median of means"
(Idea 1 was just one big mean
or \approx "mean of means")



- Z_i 's independent & the probability each Z_i out of "good range" is at most $\frac{8}{9}$
- Intuition: For large k' , vast majority of Z_i 's should be in the good range
- The median not in "good range" only if half of estimates above the range or below it