

# Multiplicative vs. Additive Error

For point queries, we wanted to return, for any  $i \in [n]$ , an estimate  $\hat{f}_i$  of  $f_i$  such that  $|f_i - \hat{f}_i| \leq \epsilon \cdot m$ . This is called an additive approximation with error  $\epsilon \cdot m$ .

Often, we will be interested in multiplicative error.

Definition: Let  $A(\sigma)$  denote the output of a randomized streaming algorithm  $A$  on input  $\sigma$ . Let  $\phi$  be the function it is supposed to compute. We say  $A$  outputs a  $(\epsilon, \delta)$ -~~additive~~ multiplicative approximation to  $\phi$  if we have for all streams  $\sigma$ :

$$\Pr \left[ \left| \frac{A(\sigma)}{\phi(\sigma)} - 1 \right| \geq \epsilon \right] \leq \delta$$

We also say a number  $N$  is a  $(\pm \epsilon)$ -multiplicative approximation to  $\phi(\sigma)$  if  $|N - \phi(\sigma)| \leq \epsilon \cdot |\phi(\sigma)|$ .

Note: Multiplicative error means if the "right answer" is small, the additive error must also be small. ~~the squared additive~~

~~the squared additive~~ In particular, if  $\phi(\sigma) = 0$ , then the answer must be exact.  
Sometimes this is impossible to achieve without storing the whole data stream (e.g. for point queries).

# AMS Sampling for Frequency-Based Functions in ~~the~~ the Vanilla Streaming model

Problem: Estimate  $\sum_{i=1}^n g(f_i)$ , where  $g$  is any function with  $g(0) = 0$ . (Want an  $(\epsilon, \delta)$ -multiplicative approximation).

Outline: Give an estimator that equals the right answer in expectation. Bound its variance. Average a bunch of copies of it to drive down the variance. Apply Chebyshev's inequality to conclude it's within a  $(1 \pm \epsilon)$ -factor of the right answer with probability  $3/4$ . ~~DO~~ this a bunch of times and take the median to drive failure prob below  $\delta$ .

Basic Estimator (unbiased, but high variance): from stream  $\langle a_1, \dots, a_m \rangle$ , sample ~~a value~~ <sup>a value  $j$</sup>  at random <sup>from  $\{1, \dots, m\}$</sup>  and compute

$$r = |\{l \geq j : a_l = a_j\}| = \# \text{ of times } a_j \text{ occurs in stream after update } j.$$

$$\text{Output } X := m \cdot (g(r) - g(r-1)).$$

$$(1, \epsilon, \delta): \mathbb{E}[X] = \sum_{i=1}^n g(f_i).$$

$$\text{Proof: } \mathbb{E}[X] = \sum_{i=1}^n \Pr[a_j = i] \cdot \mathbb{E}[X | a_j = i]$$

To see this, note:

$$\begin{aligned} &= \sum_{i=1}^n \frac{f_i}{m} \cdot \sum_{r=1}^{f_i} \frac{m \cdot (g(r) - g(r-1))}{f_i} \\ &= \sum_{i=1}^n \Pr[a_j = i] \cdot \mathbb{E}[X | a_j = i] \\ &= \sum_{i=1}^n \Pr[a_j = i] \cdot \sum_z \Pr[X=z | a_j = i] \\ &= \sum_{i=1}^n \Pr[a_j = i] \cdot \sum_z \frac{\Pr[X=z \text{ and } a_j = i]}{\Pr[a_j = i]} = \sum_{i=1}^n \sum_z z \cdot \Pr[X=z \text{ and } a_j = i] = \sum_z z \cdot \sum_{i=1}^n \Pr[X=z \text{ and } a_j = i] \\ &= \sum_z z \cdot \Pr[X=z] = \mathbb{E}[X] \end{aligned}$$

*used that  $g(0) = 0$*

Let us now look at AMS sampling specifically, in the context of  $g(r) = r^k$ . That is, we want to compute  $F_k = \sum_{i=1}^n f_i^k$ .

$F_k$  is called the  $k$ 'th frequency moment of the stream.

E.g. if the frequency vector is interpreted as the empirical distribution of an underlying distribution  $\tau$  (i.e. all stream updates are interpreted as a random draw from  $\tau$ ), then  $F_2$  is proportional to the sample variance i.e.  $\sum_{i=1}^n f_i^2 - \left(\sum_{i=1}^n f_i\right)^2 = F_2 - m^2$ .

$F_2$  also has applications to data bases (self-join size)

Claim 2:  $\text{Var}[X] \leq K \cdot n^{1/k} \cdot F_k$ . Intuitively, this means a standard deviation for basic estimator is  $\sqrt{K} \cdot n^{1/k}$  times larger than what we want (which is about  $F_k$ ). So will need to average about  $K \cdot n^{1/k}$  copies.

Before proving Claim 2, let us explain why it is enough to give an  $(\epsilon, \delta)$ -multiplicative approximation with

space cost  $O\left(\frac{1}{\epsilon^2} \cdot \log\left(\frac{1}{\delta}\right) \cdot K \cdot n^{1-1/k} \cdot (\log m \log n)\right)$

For constant  $K$ , and using  $\tilde{O}$  notation to hide factors polylogarithmic in  $m, n$ , and  $\frac{1}{\delta}$ , this is  $\tilde{O}\left(\frac{n^{1-1/k}}{\epsilon^2}\right)$ . Later in the course we will

see an algorithm achieving optimal space  $\tilde{O}\left(\frac{n^{1-2/k}}{\epsilon^2}\right)$ .

• Final estimator: First, drive down variance further by averaging several independent copies of the Basic Estimator. Call this the Mean Estimator. Second, output the median of several independent copies of the Mean Estimator.

• Analysis: Chebyshev's inequality implies the Mean Estimator is a  $(1 \pm \epsilon)$ -approximation. Chernoff Bounds imply the Median-of-Means estimator is a  $(1 \pm \epsilon)$ -approximation.

more precisely, use the following version of the Chernoff bound for additive error.

Let  $e_1, \dots, e_t$  be <sup>i.i.d. random</sup> variables taking values in  $[0, 1]$ , and let  $\mu = \sum_{i=1}^t \mathbb{E}[e_i]$

and  $e = \sum_{i=1}^t e_i$ . Then  $\Pr[|X - \mu| \geq \lambda] \leq 2 \exp\left(-\frac{2\lambda^2}{t}\right)$

So when  $\lambda = \frac{t}{4}$ , we get the probability is at most  $2 \exp\left(-\frac{2\left(\frac{t}{4}\right)^2}{t}\right)$   
 $= \exp(-8t)$ . For  $t = \frac{1}{8} \log\left(\frac{1}{\delta}\right)$ , this is  $\delta$ .

Claim 3: There is a universal constant  $c$  such that the following holds. Let  $X$  be an unbiased estimator for a real quantity  $Q$ . Let

$\{X_{i,j}\}_{i \in [t], j \in [d]}$  be a collection of i.i.d. variables each distributed identically to  $X$ , where  $t = c \cdot \log(\frac{1}{\delta})$  and  $d = \frac{4 \cdot \text{Var}[X]}{\epsilon^2 \cdot \mathbb{E}[X]}$ . (can't get the Fr. in Var)

Let  $Z = \text{median}_{i \in [t]} \left( \frac{1}{d} \sum_{j=1}^d X_{i,j} \right)$ . Then  $\Pr[|Z - Q| \geq \epsilon Q] \leq \delta$ . by claim, this is at most  $\frac{4 \cdot \text{Var}[X]}{\epsilon^2}$

Proof: By linearity of expectation, For each  $i \in [t]$ , let  $Y_i := \frac{1}{d} \sum_{j=1}^d X_{i,j}$ .

$$\mathbb{E}[Y_i] = \frac{1}{d} \sum_{j=1}^d \mathbb{E}[X_{i,j}] = \frac{1}{d} \sum_{j=1}^d Q = Q.$$

Since the variables  $X_{i,j}$  are independent,

$$\text{Var}[Y_i] = \frac{1}{d^2} \sum_{j=1}^d \text{Var}[X_{i,j}] = \frac{\text{Var}[X]}{d}.$$

By Chebyshev,  $\Pr[|Y_i - Q| \geq \epsilon Q] \leq \frac{\text{Var}[Y_i]}{(\epsilon \cdot Q)^2} = \frac{\text{Var}[X]}{d \cdot \epsilon^2 \cdot \mathbb{E}[X]}$

to see this, note that

$$\epsilon \cdot Q = \frac{\epsilon \cdot Q \cdot \sigma(Y_i)}{\sigma(Y_i)}$$

so Chebyshev says the LHS is at most the boxed quantity squared.

where  $\sigma$  is the

Hence, Chernoff bounds tell us that the probability the median of the  $Y_i$ 's has error more than  $\epsilon \cdot Q$  is at most  $\delta$ .

(Apply Chernoff bounds to the variables  $e_1, \dots, e_t$  where  $e_i = 1$  if  $|Y_i - Q| \geq \epsilon Q$ , and 0 otherwise.)

Proof of (a.m)<sup>2</sup>:

$$\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \leq \mathbb{E}[X^2]$$

$$= \sum_{j=1}^n \frac{f_j}{m} \sum_{r=1}^{f_j} \frac{1}{f_j} \cdot m^2 (r^k - (r-1)^k)^2 = m \sum_{j=1}^n \sum_{r=1}^{f_j} (r^k - (r-1)^k)^2$$

By the mean value theorem, there exists a  $\tau(z) \in [z-1, z]$  such that  $z^k - (z-1)^k = k \cdot \tau(z)^{k-1} \leq k \cdot z^{k-1}$ , where the last step requires  $k \geq 1$ .

$$\text{Hence } \hookrightarrow \leq m \sum_{j=1}^n \sum_{r=1}^{f_j} k \cdot r^{k-1} \cdot (r^k - (r-1)^k)$$

$$\leq m \sum_{j=1}^n k \cdot f_j^{k-1} \sum_{r=1}^{f_j} (r^k - (r-1)^k)$$

$$= m \sum_{j=1}^n k f_j^{k-1} f_j^k = k \cdot m \cdot \sum_{j=1}^n f_j^{2k-1}$$

$$= k \cdot \left( \sum_{j=1}^n f_j \right) \cdot \sum_{j=1}^n (f_j^{2k-1})$$

Lemma:  $\hookrightarrow \leq n^{1-1/k} \cdot F_k^2$

Proof: Let  $f_{\max} = \max_i f_i$ . Then  $f_{\max}^{k-1} = (f_{\max}^k)^{\frac{k-1}{k}} \leq \left( \sum_{i=1}^n f_i^k \right)^{\frac{k-1}{k}}$

By convexity of the function  $x \mapsto x^k$ , we have

$$\frac{1}{n} \sum_i f_i \leq \left( \frac{1}{n} \sum_i f_i^k \right)^{1/k} \quad \text{i.e.} \quad \left( \sum_i f_i \right) \leq n^{1-1/k} \cdot \sum_i f_i^k$$

biggest gap between  $\sum_i f_i$  &  $\sum_i f_i^k$  is all  $f_i$  equal  $\frac{1}{n}$   
 $\left( \sum_{i=1}^n 1 \right) \leq n^{1-1/k} \cdot \sum_{i=1}^n 1$

Hence 
$$\left( \sum_i f_i \right) \left( \sum_i f_i^{2k-1} \right) \leq \left( \sum_i f_i \right) \cdot \left( f_{\max}^{k-1} \sum_i f_i^k \right)$$

$$\leq \left( \sum_i f_i \right) \cdot \left( \sum_i f_i^k \right)^{\frac{k-1}{k}} \cdot \left( \sum_i f_i^k \right)$$

$$\leq h^{1-1/k} \cdot \left( \sum_i f_i^k \right)$$

$$= h^{1-1/k} \cdot \left( \sum_i f_i^k \right)^2 = h^{1-1/k} \cdot F_k^2$$