

# Mining Data Streams-Estimating Frequency Moment

Barna Saha

October 26, 2017

# Frequency Moment

- ▶ Computing “moments” involves distribution of frequencies of different elements in the stream.

# Frequency Moment

- ▶ Computing “moments” involves distribution of frequencies of different elements in the stream.
- ▶ Let  $f_i$  be the number of occurrences of the  $i$ th element for any  $i \in [1, n]$ , then the  $k$ th frequency moment is  $F_k = \sum_i f_i^k$

# Frequency Moment

- ▶ The 0th moment is the sum of 1 for each  $f_i > 0$ . Hence it counts the number of distinct items.

# Frequency Moment

- ▶ The 0th moment is the sum of 1 for each  $f_i > 0$ . Hence it counts the number of distinct items.
- ▶ The 1st moment is the sum of the  $f_i$ s which must be the length of the stream. This is easy to calculate.

# Frequency Moment

- ▶ The 0th moment is the sum of 1 for each  $f_i > 0$ . Hence it counts the number of distinct items.
- ▶ The 1st moment is the sum of the  $f_i$ s which must be the length of the stream. This is easy to calculate.
- ▶ The 2nd moment is the sum of the squares of the  $f_i$ 's. It is sometimes called the *surprise number* as it measures the unevenness of the distribution of elements.
  - ▶ Suppose we have a stream of length 100.
  - ▶ Scenario 1: There are 10 elements each with frequency 10.  
 $F_2 = 10 * 10^2 = 1000$
  - ▶ Scenario 2: There are 10 elements, 1st item has frequency 91, and rest have each frequency 1.  $F_2 = 91^2 + 9 * 1^2 = 8290$ .

# Computing $F_2$ in Small Space

- ▶ Alon-Matias-Szegedy: Linear Sketching

## Linear Sketch for $F_2$

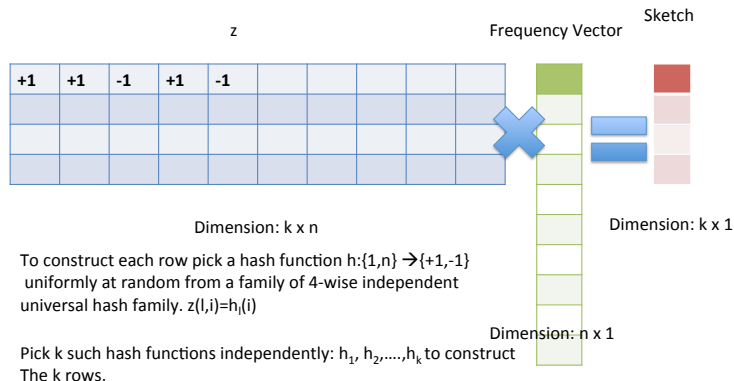
- ▶ **Problem** Given a stream  $A_1, A_2, \dots, A_m$  where elements are coming from the universe  $[1, n]$  estimate  $F_2 = \sum_{i=1}^n f_i^2$  in “small space”.
- ▶ **Output** Return an estimate  $\hat{F}_2$  such that

$$\Pr\left(F_2(1 - \epsilon) \leq \hat{F}_2 \leq (1 + \epsilon)F_2\right) \geq (1 - \delta)$$

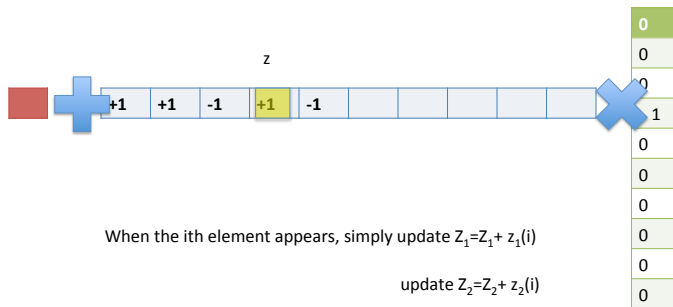
where  $\epsilon > 0$  and  $\delta > 0$  are respectively the error and confidence parameters.



# Linear Sketch for $F_2$



# Linear Sketch for $F_2$



When the  $i$ th element appears, simply update  $Z_1 = Z_1 + z_1(i)$

update  $Z_2 = Z_2 + z_2(i)$

update  $Z_3 = Z_3 + z_3(i)$

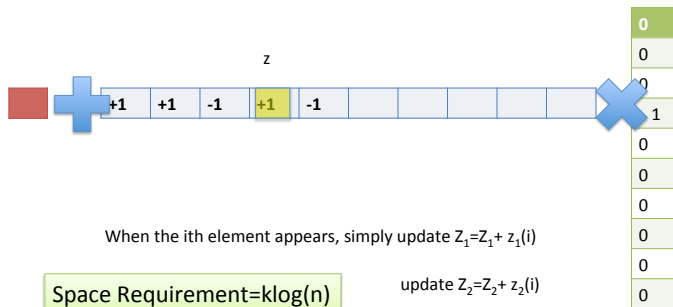
$\vdots$

$\vdots$

$\vdots$

update  $Z_k = Z_k + z_k(i)$

# Linear Sketch for $F_2$



When the  $i$ th element appears, simply update  $Z_1 = Z_1 + z_1(i)$

Space Requirement =  $k \log(n)$

Estimate =  $(Z_1^2 + Z_2^2 + \dots + Z_k^2) / k$

update  $Z_2 = Z_2 + z_2(i)$

update  $Z_3 = Z_3 + z_3(i)$

$\vdots$

$\vdots$

update  $Z_k = Z_k + z_k(i)$

Estimate:  $\hat{F}_2 = \frac{1}{k} \sum_{i=1}^k Z_i^2$

- ▶ Why is this a good estimate?

Estimate:  $\hat{F}_2 = \frac{1}{k} \sum_{i=1}^k Z_i^2$

- ▶ Why is this a good estimate?
- ▶ Show  $E[\hat{F}_2] = F_2$ .

Estimate:  $\hat{F}_2 = \frac{1}{k} \sum_{i=1}^k Z_i^2$

- ▶ Why is this a good estimate?
- ▶ Show  $E[\hat{F}_2] = F_2$ .
- ▶ Show  $\text{Var}[\hat{F}_2] \leq \frac{2F_2^2}{k}$ .

$$\text{Estimate: } \hat{F}_2 = \frac{1}{k} \sum_{i=1}^k Z_i^2$$

- ▶ Why is this a good estimate?
- ▶ Show  $E[\hat{F}_2] = F_2$ .
- ▶ Show  $\text{Var}[\hat{F}_2] \leq \frac{2F_2^2}{k}$ .
- ▶ Apply Chebyshev.

$$\text{Prob} \left( |\hat{F}_2 - F_2| > \epsilon F_2 \right) \leq \frac{\text{Var}(\hat{F}_2)}{\epsilon^2 F_2^2}$$

$$\text{Estimate: } \hat{F}_2 = \frac{1}{k} \sum_{i=1}^k Z_i^2$$

- ▶ Why is this a good estimate?
- ▶ Show  $E[\hat{F}_2] = F_2$ .
- ▶ Show  $\text{Var}[\hat{F}_2] \leq \frac{2F_2^2}{k}$ .
- ▶ Apply Chebyshev.

$$\text{Prob} \left( |\hat{F}_2 - F_2| > \epsilon F_2 \right) \leq \frac{\text{Var}(\hat{F}_2)}{\epsilon^2 F_2^2}$$

- ▶ Take  $k = \frac{16}{\epsilon^2}$ .  $\text{Prob} \left( |\hat{F}_2 - F_2| > \epsilon F_2 \right) \leq \frac{1}{8}$



$$\text{Estimate: } \hat{F}_2 = \frac{1}{k} \sum_{i=1}^k Z_i^2$$

- ▶ Why is this a good estimate?
- ▶ Show  $E[\hat{F}_2] = F_2$ .
- ▶ Show  $\text{Var}[\hat{F}_2] \leq \frac{2F_2^2}{k}$ .
- ▶ Apply Chebyshev.

$$\text{Prob} \left( |\hat{F}_2 - F_2| > \epsilon F_2 \right) \leq \frac{\text{Var}(\hat{F}_2)}{\epsilon^2 F_2^2}$$

- ▶ Take  $k = \frac{16}{\epsilon^2}$ .  $\text{Prob} \left( |\hat{F}_2 - F_2| > \epsilon F_2 \right) \leq \frac{1}{8}$

▶

$$\text{Prob} \left( F_2(1 - \epsilon) \leq \hat{F}_2 \leq (1 + \epsilon)F_2 \right) \geq \frac{7}{8}$$

## Expectation of $Z_s^2$

$$Z_s \sim Z, s = 1, 2, \dots, k$$

$$\blacktriangleright Z = \sum_{i=1}^n f_i z(i), \quad Z^2 = \sum_{i,j \in [1,n]} f_i f_j z_i z_j$$

## Expectation of $Z_s^2$

$Z_s \sim Z, s = 1, 2, \dots, k$

- ▶  $Z = \sum_{i=1}^n f_i z(i), Z^2 = \sum_{i,j \in [1,n]} f_i f_j z_i z_j$
- ▶  $E[Z^2] = \sum_{i,j \in [1,n]} E[f_i f_j z_i z_j] = \sum_i E[f_i^2 z_i^2] = \sum_i f_i^2 = F_2$

since  $E[z_i z_j] = 0$  if  $i \neq j$  and  $E[z_i^2] = 1$ .

## Expectation of $Z_s^2$

$Z_s \sim Z, s = 1, 2, \dots, k$

- ▶  $Z = \sum_{i=1}^n f_i z(i), Z^2 = \sum_{i,j \in [1,n]} f_i f_j z_i z_j$
- ▶  $E[Z^2] = \sum_{i,j \in [1,n]} E[f_i f_j z_i z_j] = \sum_i E[f_i^2 z_i^2] = \sum_i f_i^2 = F_2$

since  $E[z_i z_j] = 0$  if  $i \neq j$  and  $E[z_i^2] = 1$ .

▶

$$E[\hat{F}_2] = \frac{1}{k} \sum_{s=1}^k E[Z_s^2] = F_2$$

## Variance of $Z_s^2$

▶  $\text{Var}(Z^2) = E[Z^4] - (E[Z^2])^2$

## Variance of $Z_s^2$

▶  $\text{Var}(Z^2) = E[Z^4] - (E[Z^2])^2$



$$\begin{aligned} E[Z^4] &= \sum_i f_i^4 E[z_i^4] + \sum_{i,j:i < j} \binom{4}{2} f_i^2 f_j^2 E[z_i^2 z_j^2] \\ &= \sum_i f_i^4 + 6 \sum_{i,j:i < j} f_i^2 f_j^2 \end{aligned}$$

since  $E[z_i z_j z_k z_l] = 0$  if  $i < j < k < l$  or 3 of the terms are equal.

## Variance of $Z_s^2$

▶  $\text{Var}(Z^2) = E[Z^4] - (E[Z^2])^2$



$$\begin{aligned} E[Z^4] &= \sum_i f_i^4 E[z_i^4] + \sum_{i,j:i < j} \binom{4}{2} f_i^2 f_j^2 E[z_i^2 z_j^2] \\ &= \sum_i f_i^4 + 6 \sum_{i,j:i < j} f_i^2 f_j^2 \end{aligned}$$

since  $E[z_i z_j z_k z_l] = 0$  if  $i < j < k < l$  or 3 of the terms are equal.



$$(E[Z^2])^2 = \left( \sum_i f_i^2 \right)^2 = \sum_i f_i^4 + 2 \sum_{i,j:i < j} f_i^2 f_j^2$$

## Variance of $Z_s^2$

▶  $\text{Var}(Z^2) = E[Z^4] - (E[Z^2])^2$



$$\begin{aligned} E[Z^4] &= \sum_i f_i^4 E[z_i^4] + \sum_{i,j:i < j} \binom{4}{2} f_i^2 f_j^2 E[z_i^2 z_j^2] \\ &= \sum_i f_i^4 + 6 \sum_{i,j:i < j} f_i^2 f_j^2 \end{aligned}$$

since  $E[z_i z_j z_k z_l] = 0$  if  $i < j < k < l$  or 3 of the terms are equal.



$$(E[Z^2])^2 = \left( \sum_i f_i^2 \right)^2 = \sum_i f_i^4 + 2 \sum_{i,j:i < j} f_i^2 f_j^2$$



$$\text{Var}(Z^2) = 4 \sum_{i,j:i < j} f_i^2 f_j^2 \leq 2F_2^2$$



## Variance of $\hat{F}_2$

$$\begin{aligned}\text{Var}(\hat{F}_2) &= \text{Var}\left(\frac{1}{k} \sum_{s=1}^k Z_s^2\right) \\ &= \frac{1}{k^2} \text{Var}\left(\sum_{s=1}^k Z_s^2\right) \text{ since } \text{Var}(aX) = a^2 \text{Var}(X) \text{ for any constant } a. \\ &= \frac{1}{k^2} \sum_{s=1}^k \text{Var}(Z_s^2) \leq \frac{1}{k^2} 2kF_2^2 = \frac{2F_2^2}{k}\end{aligned}$$

# Boosting Confidence by Median

- ▶ We have

$$\text{Prob} \left( F_2(1 - \epsilon) \leq \hat{F}_2 \leq (1 + \epsilon)F_2 \right) \geq \frac{7}{8}$$

- ▶ We want

$$\text{Prob} \left( F_2(1 - \epsilon) \leq \hat{F}_2 \leq (1 + \epsilon)F_2 \right) \geq 1 - \delta$$

# Boosting Confidence by Median

- ▶ We have

$$\text{Prob} \left( F_2(1 - \epsilon) \leq \hat{F}_2 \leq (1 + \epsilon)F_2 \right) \geq \frac{7}{8}$$

- ▶ We want

$$\text{Prob} \left( F_2(1 - \epsilon) \leq \hat{F}_2 \leq (1 + \epsilon)F_2 \right) \geq 1 - \delta$$

- ▶ Take  $t$  independent estimates

$$H_1 = \hat{F}_2^1, H_2 = \hat{F}_2^2, \dots, H_t = \hat{F}_2^t$$

# Boosting Confidence by Median

- ▶ We have

$$\text{Prob} \left( F_2(1 - \epsilon) \leq \hat{F}_2 \leq (1 + \epsilon)F_2 \right) \geq \frac{7}{8}$$

- ▶ We want

$$\text{Prob} \left( F_2(1 - \epsilon) \leq \hat{F}_2 \leq (1 + \epsilon)F_2 \right) \geq 1 - \delta$$

- ▶ Take  $t$  independent estimates

$$H_1 = \hat{F}_2^1, H_2 = \hat{F}_2^2, \dots, H_t = \hat{F}_2^t$$

- ▶ Return the median of  $H_1, H_2, \dots, H_t$ .

## Boosting by Median

- ▶ Suppose there is an Algorithm that returns an estimate  $\hat{F}$  of a true estimate  $F$  such that  $|\hat{F} - F|$  is small with probability  $\frac{7}{8}$ .
- ▶ How can we design an algorithm that will return an estimate  $G$  of  $F$  such that  $|G - F|$  is small with probability  $99/100$ ?  
(In general  $1 - \delta$ )

## Boosting by Median

- ▶ Suppose there is an Algorithm that returns an estimate  $\hat{F}$  of a true estimate  $F$  such that  $|\hat{F} - F|$  is small with probability  $\frac{7}{8}$ .
- ▶ How can we design an algorithm that will return an estimate  $G$  of  $F$  such that  $|G - F|$  is small with probability  $99/100$ ? (In general  $1 - \delta$ )
- ▶ Run  $s = 2 \log \frac{2}{\delta} + 1$  independent copies of the Algorithm to obtain estimates  $\hat{F}^1, \hat{F}^2, \dots, \hat{F}^s$ . Set  $G = \text{median}_{i=1}^s \hat{F}^i$ .

# Boosting by Median

- ▶ What is the probability that the median is a bad estimate?

# Boosting by Median

- ▶ What is the probability that the median is a bad estimate?
- ▶ Either all  $\lfloor \frac{s}{2} \rfloor$  copies with estimate below  $G$  are bad or,  $\lfloor \frac{s}{2} \rfloor$  copies with estimate above  $G$  are bad. That is there are  $\log \frac{2}{\delta}$  copies that are at least bad for  $G$  to be a bad estimate.



# Boosting by Median

- ▶ What is the probability that the median is a bad estimate?
- ▶ Either all  $\lfloor \frac{s}{2} \rfloor$  copies with estimate below  $G$  are bad or,  $\lfloor \frac{s}{2} \rfloor$  copies with estimate above  $G$  are bad. That is there are  $\log \frac{2}{\delta}$  copies that are at least bad for  $G$  to be a bad estimate.
- ▶ Show that the probability of Median to be bad is  $\leq \delta$

# Frequency Moment

- ▶ For  $k > 2$ , the best bound known is  $\tilde{O}(n^{1-\frac{2}{k}} \log \frac{1}{\delta})$  barring  $\text{poly}(\frac{1}{\epsilon})$  factor. There is an almost matching lower bound of  $\Omega(n^{1-\frac{2}{k}})$ .
- ▶ For  $k < 2$ , the best bound known is  $\tilde{O}(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ .
- ▶ The algorithms use clever combination of sketching and hashing

# Sketching as a Versatile Tool

- ▶ Estimating entropy, quantiles, heavy hitters, fitting histograms etc.
- ▶ Applications beyond streaming: dimensionality reduction, nearest neighbors, anomaly detection, statistics over social network.
- ▶ Not only useful for small-space algorithm design, but also for fast running time, distributed processing etc.

# Sketching as a Versatile Tool

## A different linear sketch

- Instead of  $\pm 1$ , let  $r_i$  be i.i.d. random variables from  $N(0,1)$
- Consider

$$Z = \sum_i r_i x_i$$

- We still have that  $E[Z^2] = \sum_i x_i^2 = \|x\|_2^2$ , since:
  - $E[r_i] E[r_j] = 0$
  - $E[r_i^2] = \text{variance of } r_i, \text{ i.e., } 1$
- As before we maintain  $\mathbf{Z} = [Z_1 \dots Z_k]$  and define
$$Y = \|\mathbf{Z}\|_2^2 = \sum_j Z_j^2 \quad (\text{so that } E[Y] = k\|x\|_2^2)$$
- We show that there exists  $C > 0$  s.t. for small enough  $\epsilon > 0$

$$\Pr[|Y - k\|x\|_2^2| > \epsilon k\|x\|_2^2] \leq \exp(-C \epsilon^2 k)$$

Slide from Piotr Indyk's course on Streaming, Sketching and Compressed Sensing

# Sliding Window Model

- ▶ Only the last  $W$  items matter where  $W$  is the window size.

# Sliding Window Model

- ▶ Only the last  $W$  items matter where  $W$  is the window size.
- ▶ Can you extend Bloom Filter, FM sketch in this setting?

# Sliding Window Model

- ▶ Only the last  $W$  items matter where  $W$  is the window size.
- ▶ Can you extend Bloom Filter, FM sketch in this setting?
- ▶ Can you extend Count-Min sketch or linear sketching techniques in this setting?

# Decaying Window Model

- ▶ No fixed window size, but older items have less importance.



# Decaying Window Model

- ▶ No fixed window size, but older items have less importance.
- ▶ Can you extend Bloom Filter, FM sketch in this setting?

# Decaying Window Model

- ▶ No fixed window size, but older items have less importance.
- ▶ Can you extend Bloom Filter, FM sketch in this setting?
- ▶ Can you extend Count-Min sketch or linear sketching techniques in this setting?