

STREAMING ALGORITHMS

- Streaming Model of Computation
- Streaming Algorithms and DFA
- Stream: Motivation and Applications
- Synopsis: Sliding Window, Histogram, Wavelets
- Sampling from Stream: Reservoir Sampling
- Linear Sketch
- Count-Min Sketch
- AMS Sketch

IMDAD ULLAH KHAN

AMS Sketch

AMS sketch to estimate second frequency moment of a stream

- The AMS Sketch (Alon, Mathias, Szegedy, 1996)
- A sketch to estimate F_2 (paper has other algorithms for higher moments)

$$\mathcal{S} = \langle a_1, a_2, a_3, \dots, a_m \rangle \quad a_i \in [n]$$

f_i : frequency of i in \mathcal{S} $\mathbf{F} = (f_1, f_2, \dots, f_n)$

$$F_2 = \sum_{i=1}^n f_i^2$$

▷ second frequency moment

Easy to compute if we store F

▷ $O(n)$ integers

Can store $f_1 + f_2 + \dots + f_n$

▷ $O(1)$ integers

Also easy $(f_1 + f_2 + \dots + f_n)^2$

AMS sketch to estimate second frequency moment of a stream

$$F_2 = \sum_{i=1}^n f_i^2$$

Can store $f_1 + f_2 + \dots + f_n$ ▷ $O(1)$ integers

$(f_1 + f_2 + \dots + f_n)^2$ can be computed by the following algorithm

Algorithm : Compute square of sum of frequencies (S)

$X \leftarrow 0$ ▷ sketch consists of 1 integer

On input a_i

$X \leftarrow X + 1$

return X^2

$$X^2 = (f_1 + f_2 + \dots + f_n)^2$$

▷ Square of sum of frequencies, we want sum of squares of frequencies

AMS sketch

AMS sketch to estimate second frequency moment of a stream

$$F_2 = \sum_{i=1}^n f_i^2 = \underline{f_1^2 + f_2^2 + \dots + f_n^2} \quad \triangleright \text{We want this}$$

$$(f_1 + f_2 + \dots + f_n)^2 \quad \triangleright \text{Easy but overestimate}$$

$$(f_1 + f_2 + f_3 + f_4)^2 = \underline{f_1^2 + f_2^2 + f_3^2 + f_4^2} + \underbrace{2(f_1f_2 + f_1f_3 + f_2f_3 + f_1f_4 + f_2f_4 + f_3f_4)}_{\text{error}}$$

What if we randomly add/subtract frequencies

$$(f_1 - f_2 + f_3 - f_4)^2 = \underline{f_1^2 + f_2^2 + f_3^2 + f_4^2} + \underbrace{2(-f_1f_2 + f_1f_3 - f_2f_3 - f_1f_4 + f_2f_4 - f_3f_4)}_{\text{error}}$$

AMS sketch

AMS sketch to estimate second frequency moment of a stream

What if we randomly add/subtract frequencies

$$(f_1 - f_2 + f_3 - f_4)^2 = \underbrace{f_1^2 + f_2^2 + f_3^2 + f_4^2}_{\text{error}} + 2(-f_1f_2 + f_1f_3 - f_2f_3 - f_1f_4 + f_2f_4 - f_3f_4)$$

Algorithm : AMS sketch to estimate second frequency moment of \mathcal{S}

Pick a random hash function $g : [n] \rightarrow \{-1, +1\}$

$X \leftarrow 0$ ▷ sketch consists of 1 integer

On input a_i

$X \leftarrow X + g(a_i)$

return X^2

$$X = f_1g(1) + f_2g(2) + \dots + f_ng(n)$$

AMS sketch is an unbiased estimate of second frequency moment

$$X = f_1g(1)+f_2g(2)+\dots+f_ng(n) \quad X^2 = (f_1g(1)+f_2g(2)+\dots+f_ng(n))^2$$

$$\begin{aligned}\mathbb{E}[X^2] &= \mathbb{E}\left[\sum_i (f_i g(i))^2\right] + \mathbb{E}\left[\sum_{i \neq j} f_i g(i) f_j g(j)\right] \\ &= \mathbb{E}\left[\sum_i f_i^2 g(i)^2\right] + \mathbb{E}\left[\sum_{i \neq j} f_i f_j g(i) g(j)\right] \\ &= \sum_i f_i^2 \mathbb{E}[g(i)^2] + \sum_{i \neq j} f_i f_j \mathbb{E}[g(i)g(j)] = F_2\end{aligned}$$

$$\therefore \mathbb{E}[g(i)^2] = 1 \quad \text{and} \quad \mathbb{E}[g(i)g(j)] = 0 \quad \text{for } i \neq j$$

$$\mathbb{E}[X^2] = F_2$$

The variance of AMS sketch estimate for F_2 is bounded

$$X^2 = (f_1g(1) + f_2g(2) + \dots + f_n g(n))^2 \quad \mathbb{E}[X^2] = F_2$$

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

$$\mathbb{E}[X^4] = \mathbb{E}\left[\sum_i (f_i g(i))^4 + 6 \sum_{i \neq j} (f_i g(i))^2 f_j g(j)^2\right] + \dots$$

other terms: $\mathbb{E}[g(i)g(j)g(k)g(l)] = \mathbb{E}[g(i)^2g(j)g(k)] = \mathbb{E}[g(i)^3g(j)] = 0$

▷ 4-wise independence

$$\mathbb{E}[X^4] = \sum_i f_i^4 + 6 \sum_{i \neq j} f_i^2 f_j^2$$

$$\text{Var}(X^2) = \sum_i f_i^4 + 6 \sum_{i \neq j} f_i^2 f_j^2 - (\sum_i f_i^2)^2 = 4 \sum_{i \neq j} f_i^2 f_j^2 \leq 2F_2^2$$

Quality Specs of basic AMS Sketch

Algorithm : AMS sketch to estimate second frequency moment of \mathcal{S}

Pick a random hash function $g : [n] \mapsto \{-1, +1\}$

$X \leftarrow 0$

▷ sketch consists of 1 integer

On input a_i

$X \leftarrow X + g(a_i)$

return X^2

$$\mathbb{E}[X^2] = F_2$$

$$\text{Var}(X^2) \leq 2F_2^2$$

Amplifying the probability of basic AMS Sketch

- Keep $k = 8/\epsilon^2 \times \log(1/\delta)$ estimates, X_1, X_2, \dots, X_k
- Return \bar{X} : median of $\log(1/\delta)$ averages of groups of $2/\epsilon^2$ estimates

Algorithm : AMS sketch to estimate F_2 of $\mathcal{S}(\epsilon, \delta)$

Pick $k = 8/\epsilon^2 \times \log(1/\delta)$ random hash functions $g_j : [n] \rightarrow \{-1, +1\}$

$X \leftarrow \text{ZEROS}(k)$

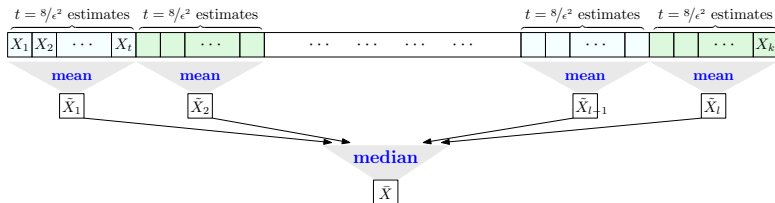
▷ sketch consists of k integer

On input a_i

for $j = 1 \rightarrow k$ **do**

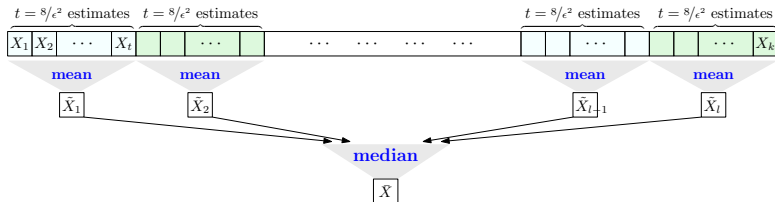
$X[j] \leftarrow X[j] + g_j(a_i)$

return \bar{X} : median of $\log(1/\delta)$ means of groups of $8/\epsilon^2$ estimates ($X[\cdot]^2$)



Amplifying the probability of basic AMS Sketch

- Keep $k = 8/\epsilon^2 \times \log(1/\delta)$ estimates, X_1, X_2, \dots, X_k
- Return \bar{X} : median of $\log(1/\delta)$ averages of groups of $2/\epsilon^2$ estimates



- $\mathbb{E}[X_j^2] = F_2$ $\text{Var}(X_j^2) \leq 2F_2^2$
- $\mathbb{E}[\tilde{X}_j] = F_2$ $\text{Var}(\tilde{X}_j) \leq \epsilon^2/4F_2^2$
- $\Pr[|\tilde{X}_j - F_2| \geq \epsilon F_2] \leq \text{Var}(\tilde{X}_j)/\epsilon^2 F_2^2 = 1/4$ \triangleright Chebyshev Inequality
- $\Pr[|\bar{X} - F_2| \geq \epsilon F_2] \leq \delta$

The last inequality uses the Chernoff bound. For \bar{X} to deviate this much from F_2 at least half of \tilde{X}_j have to deviate more than that

Linear Transformation View of AMS Sketch

Algorithm : AMS sketch to estimate F_2 of \mathcal{S}

Pick k random hash functions $g : [n] \mapsto \{-1, +1\}$

$X \leftarrow \text{ZEROS}(k)$

▷ sketch consists of 1 integer

On input a_i

for $j = 1 \rightarrow k$ **do**

$X[j] \leftarrow X[j] + g_j(a_i)$

$$\mathbf{g} = \begin{array}{|c|c|c|c|} \hline g(1) & g(2) & \dots & g(n) \\ \hline \end{array}$$

F

$$\begin{array}{|c|} \hline f_1 \\ \hline f_2 \\ \hline \vdots \\ \hline \vdots \\ \hline f_n \\ \hline \end{array}$$

$= X$

Linear Transformation View of AMS Sketch

Algorithm : AMS sketch to estimate F_2 of \mathcal{S}

Pick k random hash functions $g : [n] \mapsto \{-1, +1\}$

$X \leftarrow \text{ZEROS}(k)$

▷ sketch consists of 1 integer

On input a_i

for $j = 1 \rightarrow k$ **do**

$X[j] \leftarrow X[j] + g_j(a_i)$

$$g = \begin{array}{|c|c|c|c|c|} \hline +1 & -1 & \dots & & +1 \\ \hline \end{array}$$

F

$$\begin{array}{|c|} \hline f_1 \\ \hline f_2 \\ \hline \vdots \\ \hline \vdots \\ \hline f_n \\ \hline \end{array}$$

$= X$

Linear Transformation View of AMS Sketch

Algorithm : AMS sketch to estimate F_2 of \mathcal{S}

Pick k random hash functions $g : [n] \mapsto \{-1, +1\}$

$X \leftarrow \text{ZEROS}(k)$

▷ sketch consists of 1 integer

On input a_i

for $j = 1 \rightarrow k$ **do**

$X[j] \leftarrow X[j] + g_j(a_i)$

$$\mathbf{G} = \begin{array}{|c|c|c|c|c|} \hline +1 & -1 & \dots & & +1 \\ \hline -1 & -1 & \dots & & -1 \\ \hline \end{array}$$

$$\mathbf{F} = \begin{array}{|c|} \hline f_1 \\ \hline f_2 \\ \hline \vdots \\ \hline \vdots \\ \hline f_n \\ \hline \end{array}$$

$$\mathbf{X} = \begin{array}{|c|} \hline X_1 \\ \hline X_2 \\ \hline \end{array}$$

Linear Transformation View of AMS Sketch

Algorithm : AMS sketch to estimate F_2 of \mathcal{S}

Pick k random hash functions $g : [n] \mapsto \{-1, +1\}$

$X \leftarrow \text{ZEROS}(k)$

▷ sketch consists of 1 integer

On input a_i

for $j = 1 \rightarrow k$ **do**

$X[j] \leftarrow X[j] + g_j(a_i)$

$\mathbf{G} =$

+1	-1	...		+1
-1	-1	...		-1
⋮		...		⋮
-1	+1	...		-1

\mathbf{F}

f_1
f_2
⋮
⋮
f_n

\mathbf{X}

X_1
X_2
⋮
X_k

Linear Transformation View of AMS Sketch

$$\mathbf{G} =$$

+1	-1	...		+1
-1	-1	...		-1
\vdots		...		\vdots
-1	+1	...		-1

F

f_1
f_2
\vdots
\vdots
f_n

X

X_1
X_2
\vdots
X_k

$$\bar{X} = \frac{1}{k} \sum_{i=1}^k X_i^2 \quad \Pr [|\bar{X} - F_2| > \epsilon F_2] \leq \delta$$

With probability at least $1 - \delta$

$$(1 - \epsilon) \sum_{i=1}^n f_i^2 < \frac{1}{k} \sum_{i=1}^k X_i^2 < (1 + \epsilon) \sum_{i=1}^n f_i^2$$
$$\sqrt{(1 - \epsilon)} \|F\|_2 < \frac{1}{\sqrt{k}} \|X\|_2 < \sqrt{(1 + \epsilon)} \|F\|_2$$

AMS Sketch as a dimensionality reduction algorithm

$$\mathbf{G} =$$

+1	-1	...		+1
-1	-1	...		-1
⋮		...		⋮
-1	+1	...		-1

$$\mathbf{F}$$

f_1
f_2
⋮
⋮
⋮
f_n

$$\mathbf{X}$$

X_1
X_2
⋮
⋮
X_k

$$\sqrt{(1 - \epsilon)} \|F\|_2 < \frac{1}{\sqrt{k}} \|X\|_2 < \sqrt{(1 + \epsilon)} \|F\|_2$$

\mathbf{G} is a random linear transformation reduces the dimension of F while preserving its ℓ_2 -norm

Since G is linear it is easy to see that given $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$

$$\text{w.h.p} \quad \left\| \frac{1}{\sqrt{k}} \mathbf{G}\mathbf{u} \right\|_2 - \left\| \frac{1}{\sqrt{k}} \mathbf{G}\mathbf{v} \right\|_2 \sim \|\mathbf{u} - \mathbf{v}\|_2$$