

DATA STREAMS

- Imdadullah Khan

Data Stream

Stream Processing: Analytics on a continuous stream of data items

The goal is to draw meaningful analytics from the stream subject to
subject to

- **Single Pass:** process each item exactly once (common requirement)
- **Limited Memory:** poly-logarithmic space (in length of stream or domain)
- **Constant per item processing:** near real time
- **Arbitrary arrival order:** No assumption on distribution or order of items

Outline

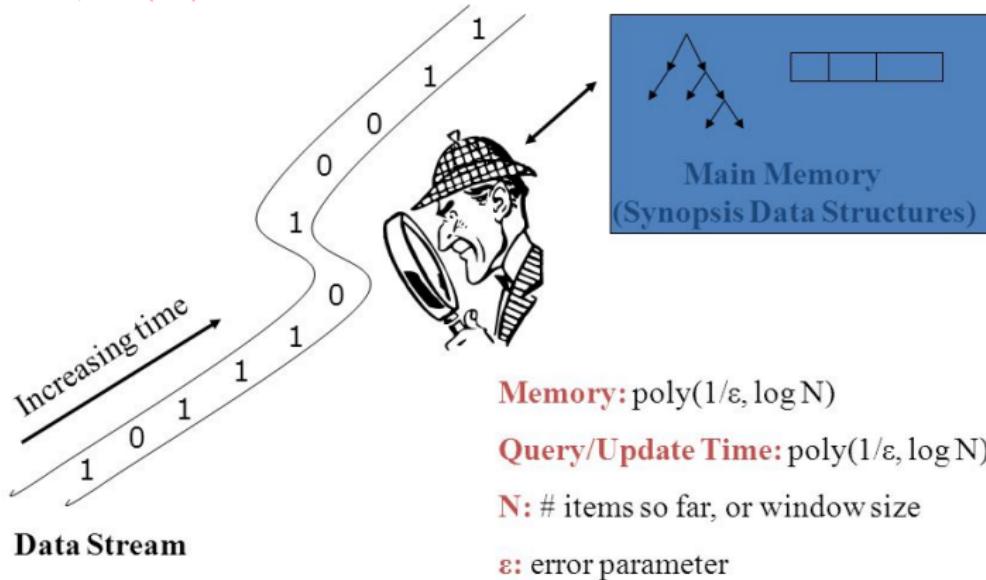
Characteristics of data streams¹

- Huge volumes of continuous data, possibly infinite
- Fast changing and requires fast, real-time response
- Data stream captures nicely our data processing needs of today
- Random access is expensive
- Single scan algorithm (can only have one look)
- Store only the summary of the data seen thus far
- Most stream data are at pretty low-level or multi-dimensional in nature, needs multi-level and multi-dimensional processing

¹Based on Han & Kamber, Data Mining Concepts & Techniques, 2nd Ed.

Stream Model of Computation

Motwani, PODS (2002)



Data Stream

Stream data is fundamentally different than traditional datasets²

Traditional Data (DBMS)	Data Stream
Persistent storage	Transient stream(s)
One-time query	Continuous query
Random access	Sequential access
Unbounded disk storage	Bounded main memory
Only current state matters	Arrival-order is critical
No real time services	Real-time requirements
Low update rate	Possibly multi-GB arrival rate (dynamic & fast)
Mixed granularity	Data at fine granularity

²R. Motwani, PODS (2002)

Data Stream Processing Model

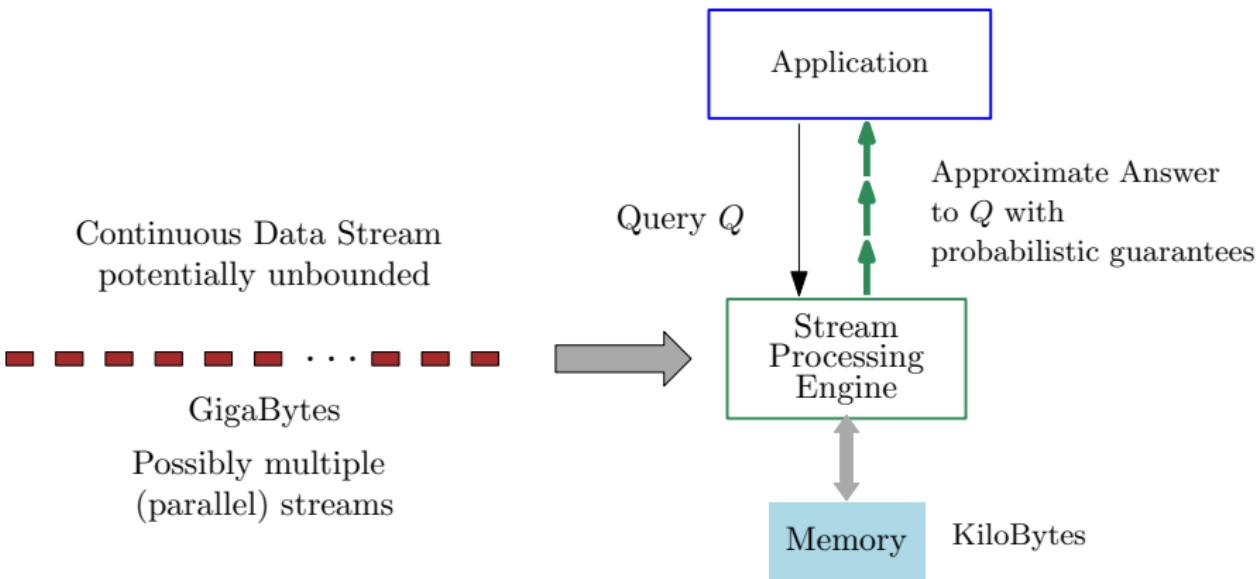
- Since streams are long (potentially unbounded) exact algorithms with limited memory are possible only for a few simple queries
- ∴ we design approximate algorithms (they often suffice)

(ϵ, δ) -approximate algorithm

- \mathcal{A} : an algorithm to compute $f(\mathcal{S})$ (a function of stream)
- $\mathcal{A}(\mathcal{S})$: output of \mathcal{A} on \mathcal{S}
- For $\epsilon > 0, 0 \leq \delta \leq 1$, \mathcal{A} is an (ϵ, δ) -approximation algorithm if

$$\Pr[\mathcal{A}(\mathcal{S}) - f(\mathcal{S}) > \epsilon f(\mathcal{S})] \leq \delta$$

Data Stream



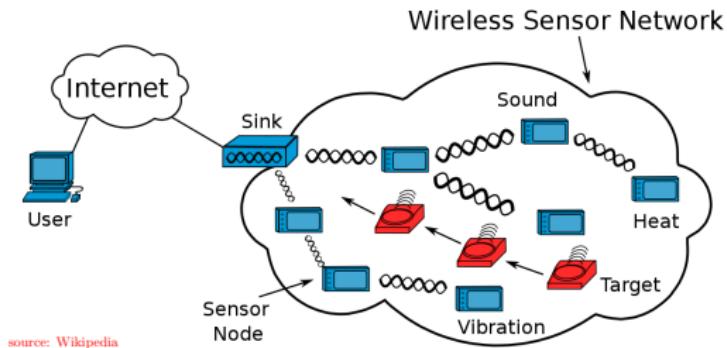
Data Stream Applications

Stream data comes in many domains and has various applications³

- Telecommunication calling records
- Business: credit card transaction flows
- Network monitoring and traffic engineering
- Financial market: stock exchange
- Engineering & industrial processes: power supply & manufacturing
- Sensor, monitoring & surveillance: video streams, RFIDs
- Security monitoring
- Web logs and Web page click streams
- Massive data sets (even saved but random access is too expensive)

³Based on Han & Kamber, Data Mining Concepts & Techniques, 2nd Ed.

Data Stream Motivation

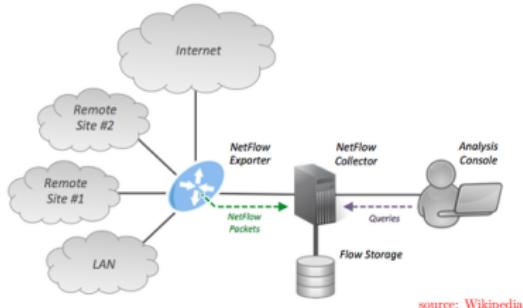


- Sensor nodes collect unlimited amount of data
- have very limited computation power and memory
- Limited battery power constrain communication of all collected data
- 1 bit transmission consumes power \sim to executing 800 instructions⁴
- Streaming algorithm deployed onto nodes are ideally suited for drawing analytics from sensed data

⁴Madden et.al. (2002)

Data Stream Applications

Network Monitoring and Management



source: Wikipedia

NetFlow: A Cisco tool for network administrators (performance metrics, security analysis, detection and forensics). For each Flow it reports (logs)

- Network Interface
- Source/Destination IP Addresses
- IP Protocol
- AT&T Processes over 567 billion flow records per day⁵ $\triangleright \sim 15$ PBytes
- Detects and characterizes approximately 500 anomalies per day

⁵ Fred Stinger (AT&T) FloCon (2017) Netflow Collection and Analysis ..

Data Stream Applications

Network Monitoring and Management

Application Areas	Quries
<ul style="list-style-type: none">■ Traffic Engineering■ Traffic Monitoring■ Volume estimation & analysis■ Load Balancing■ Efficient Resource Utilization■ (D)DOS Attack Detection■ SLA Voilation	<ul style="list-style-type: none">■ How many bytes sent b/w IP-1 and IP-2?■ How many IP addresses are active?■ Top 100 IP's by traffic volume■ Average duration of IP session?■ Meidan number of bytes in each IP session■ Find sessions that transmitted $> 1k$ bytes■ Find sessions with duration $>$ twice average■ List all IP's with a sudden spike in traffic■ List all IP involved in more than 1k sessions

Data Stream Applications

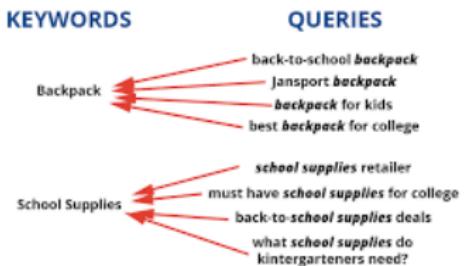
Web Click Stream Analysis: tracking and analysis of websites visits



- A stream of user clicks (and other actions) on a website
- Which webpages are getting large/small number of clicks
- Used to enhance customer experience
- Improving conversion rates
- Optimizing digital marketing
- Up-selling and cross-selling

Data Stream Applications

Search Queries Stream:



- **Search Engine:** logs of billions of queries every day
- Discover trends and patterns
- Find relevant keywords for your website
- Estimate competition scores (or difficulty) and CPC (cost per click) of keywords
- Improve Good advertisement

Data Stream Applications

Energy consumption Analysis:



- Electricity consumption data from AMI (Automatic Metering Interface)
- Find average hourly load, load surges, anomaly
- Identify faults, drops, failures

Data Stream Applications

Query Execution Plan can be optimized using a synopsis of the database

Suppose we have data of $n = 1M$ people in a database and the query

`SELECT * from Table WHERE $25 \leq age \leq 35$ and $54 \leq weight \leq 60$`

Runtime of brute force execution is $2n$ comparisons

Suppose we have the following synopsis of distribution of an attribute

Age	Freq
0 – 10	7%
11 – 20	8%
21 – 30	10%
31 – 40	12%
41 – 50	13%
51 – 60	25%
61 – 70	20%
71+	5%

First filter on Age, then on weight

Runtime: $1.22n$

Weight	Freq.
0 – 20	20%
21 – 40	25%
41 – 60	10%
61 – 80	15%
81+	30%

First filter on Weight, then on age

Runtime: $1.1n$

Stream Model of Computation

Stream $S := a_1, a_2, a_3, \dots, a_m$

▷ m may be unknown

Each $a_i \in [n]$

Goal: Compute a function of the stream S (e.g. mean, median, number of distinct elements, frequency moments..)

Subject to

- Single pass, read each element of S only once sequentially
- Per item processing time $O(1)$
- Use memory polynomial in $O(1/\epsilon, 1/\delta, \log n)$
- Return (ϵ, δ) -randomized approximate solution

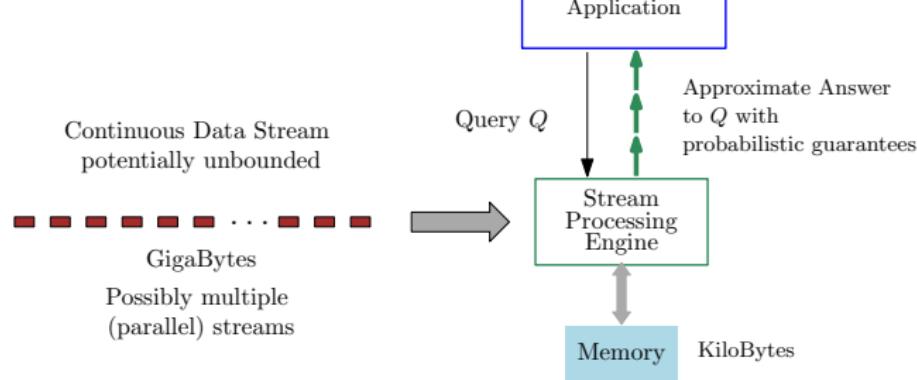
Data Stream: Synopsis

Fundamental Methodology: Keep a synopsis of the stream and answer query based on it. Update synopsis after examining each item in $O(1)$

Synopsis: Succinct summary of the stream (so far) (poly-log bits)

Families of Synopsis

- Sliding Window
- Random Sample
- Histogram
- Wavelets
- Sketch



Synopsis Based Exact Stream Computation

- Length of \mathcal{S} (m): Computed by storing a running counter
- Sum of \mathcal{S} : Computed by storing a running sum
- Mean of \mathcal{S} : Computed from sum and length of \mathcal{S}
- Variance of \mathcal{S} : Computed from sum, sum of square, and length of \mathcal{S}

$$\text{Var}(X) = E(X^2) - (E(X))^2$$

All these use $O(\log n)$ bits memory and constant time per element

Synopsis Based Exact Stream Computation

Missing Element

- $n - 1$ unique integers are streamed in from $[n]$
- Find the missing integer?
- Trivial to find it if we use n bits
- A better solution is to save sum S of the stream $\triangleright O(\log n)$ bits
- The missing integer is $n(n+1)/2 - S$
- Can do it in exactly $\log n$ bits by storing the parity sum of each bits
- The final parity sum is the missing integer

Synopsis Based Exact Stream Computation

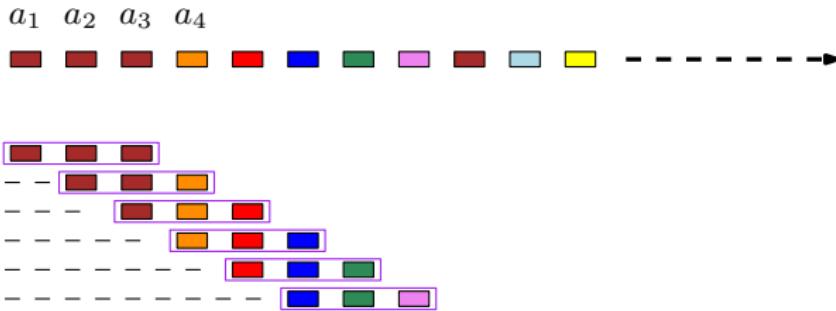
Two Missing Elements

- $n - 2$ unique integers are streamed in from $[n]$
- Find the missing integers?
- Trivial to find it if we use n bits
- Save sum of 1st and 2nd powers of stream elements $\triangleright O(\log n)$ bits
- The missing integers are solution to 2 unknowns and two equations
- Readily generalizes to k missing elements

Data Stream: Sliding Window

Synopsis: Sliding Window

- Keep the last w elements as synopsis (w is length of window)
- On input a_i ($i \geq w$), a_{i-w} expires and a_i added to window
- Can be used for queries like mean, sum, variance, count of pre-specified element(s) (e.g. non-zero, even)
- Extended to compute approximate median, and k -median

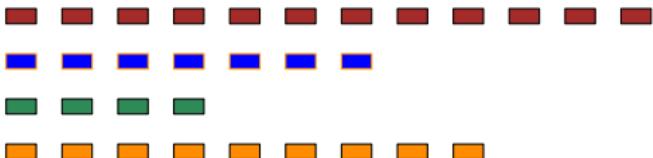


Data Stream: Random Sample

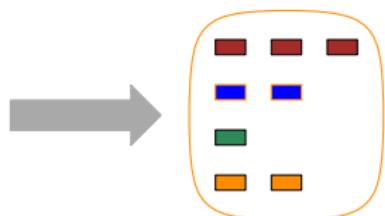
Synopsis: Random Sample

- Keep a “**representative**” subset of the stream
- Approximately compute query answer on sample (with appropriate scaling etc.)

Stream elements in an arbitrary order



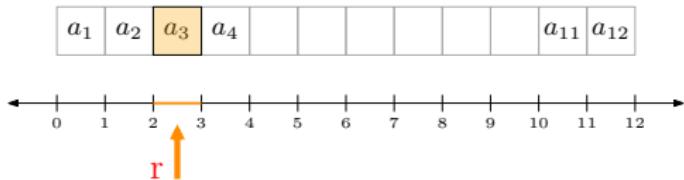
Random Sample



Data Stream: Random Sample

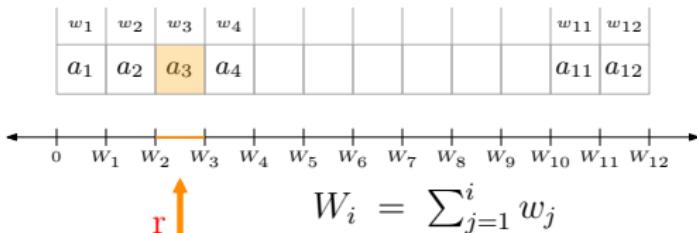
Sample a random element from array A of length n $\triangleright A[i]$ with prob $1/n$

- Generate a random number $r \in [0, n]$ $\triangleright r \leftarrow \text{RAND}() \times n$
- Return $A[\lceil r \rceil]$



Sample random element (by weight) from array A $\triangleright A[i]$ with prob. w_i/W

- Generate a random number $r \in [0, \sum_{j=1}^n w_j]$ $\triangleright r \leftarrow \text{RAND}() \times W_n$
- Return $A[i]$ if $W_{i-1} \leq r < W_i$



Data Stream: Random Sample

Sample a random element from the stream S

▷ a_i with prob. $1/m$

- If m is known, use algorithm for sampling from array. For unknown m

Algorithm : Reservoir Sampling (\mathcal{S})

$R \leftarrow a_1$ ▷ R (reservoir) maintains the sample

for $i \geq 2$ do

Pick a_i with probability $1/i$

Replace with current element in R

Prob. that a_i is in the sample R_m (m : stream length or query time)

$$= \underbrace{\Pr \text{ that } a_i \text{ was selected at time } i}_{\frac{1}{i}} \times \underbrace{\Pr \text{ that } a_i \text{ survived in } R \text{ until time } m}_{\prod_{j=i+1}^m \left(1 - \frac{1}{j}\right)}$$

$$= \cancel{\frac{1}{i}} \times \cancel{\frac{i}{i+1}} \times \cancel{\frac{i+1}{i+2}} \times \cancel{\frac{i+2}{i+3}} \times \dots \times \cancel{\frac{m-2}{m-1}} \times \cancel{\frac{m-1}{m}} = \frac{1}{m}$$

Data Stream: Random Sample

Sample k random elements from the stream S

▷ a_i with prob. k/m

Algorithm : Reservoir Sampling (S, k)

$R \leftarrow a_1, a_2, \dots, a_k$ ▷ R (reservoir) maintains the sample

for $i \geq k + 1$ **do**

- Pick a_i with probability k/i
- If a_i is picked, replace with it a randomly chosen element in R

Prob. that a_i is in the sample R_m (m : stream length or query time)

$$= \underbrace{\Pr \text{ that } a_i \text{ was selected at time } i}_{\frac{k}{i}} \times \underbrace{\Pr \text{ that } a_i \text{ survived in } R \text{ until time } m}_{\prod_{j=i+1}^m \left(1 - \left(\frac{k}{j} \times \frac{1}{k}\right)\right)}$$

$$= \cancel{\frac{k}{i}} \times \cancel{\frac{i}{i+1}} \times \cancel{\frac{i+1}{i+2}} \times \cancel{\frac{i+2}{i+3}} \times \dots \times \cancel{\frac{m-2}{m-1}} \times \cancel{\frac{m-1}{m}} = \frac{k}{m}$$

Data Stream: Histogram and Wavelets

Synopsis: Histogram

- The synopsis is some summary statistics (e.g. frequency, mean) of groups (subsets, buckets) in streams values
 - Equi-width histogram
 - Equidepth histogram
 - V-optimal histogram
 - Multi-dimensional histogram

Synopsis: Wavelets

- Essentially histograms of features (coefficients) in the frequency domain representation of the stream

Data Stream: Linear Sketch

- **Sample** is a general purpose synopsis
- Process sample only – no advantage from observing the whole stream
- Sketches are specific to a particular purpose (query)
- **Sketches (also histograms and wavelets)** take advantage from the fact the processor see the whole stream (though can't remember all)

A **linear sketch** interprets the stream as defining the frequency vector

Often we are interested in functions of the frequency vector from a stream

$$\mathcal{S} : a_1, a_2, a_3, a_4, \dots, a_m$$
$$a_i \in [n]$$

$$\mathbf{F} : \begin{array}{|c|c|c|c|c|c|c|c|c|c|c|c|} \hline 1 & 2 & 3 & & & & & & & & n \\ \hline f_1 & f_2 & f_3 & & \dots & & \dots & & & & f_n \\ \hline \end{array}$$
$$f_j = |\{a_i \in \mathcal{S} : a_i = j\}| \quad (\text{frequency of } j \text{ in } \mathcal{S})$$

$$\mathcal{S} : 2, 5, 6, 7, 8, 2, 1, 2, 7, 5, 5, 4, 2, 8, 8, 9, 5, 6, 4, 4, 2, 5, 5$$

$$\mathbf{F} : \begin{array}{|c|c|c|c|c|c|c|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \hline 1 & 5 & 0 & 3 & 6 & 2 & 2 & 3 & 1 \\ \hline \end{array}$$

Data Stream: Linear Sketch

A [linear sketch](#) interprets the stream as defining the frequency vector

Often we are interested in functions of the frequency vector from a stream

$$\mathcal{S} : a_1, a_2, a_3, a_4, \dots, a_m \quad \quad \mathbf{F} : \begin{array}{|c|c|c|c|c|c|} \hline & 1 & 2 & 3 & & n \\ \hline f_1 & f_2 & f_3 & & \dots & \dots \\ \hline \end{array} f_n$$

$f_j = |\{a_i \in \mathcal{S} : a_i = j\}|$ (frequency of j in \mathcal{S})

$$\mathcal{S} : \quad 2, 5, 6, 7, 8, 2, 1, 2, 7, 5, 5, 4, 2, 8, 8, 9, 5, 6, 4, 4, 2, 5, 5$$

$$\mathbf{F} : \begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \hline 1 & 5 & 0 & 3 & 6 & 2 & 2 & 3 & 1 \\ \hline \end{array}$$

Stream: Frequency Moments

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

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

$$F_0 := \sum_{i=1}^n f_i^0 \quad \triangleright \text{number of distinct elements}$$

$$F_1 := \sum_{i=1}^n f_i \quad \triangleright \text{length of stream, } m$$

$$F_2 := \sum_{i=1}^n f_i^2 \quad \triangleright \text{second frequency moment}$$

Data Stream: Linear Sketch

Synopsis: Linear Sketches

Linear sketch is a synopsis that can be computed as a linear transform of \mathbf{F}

- Best suited for data streams, highly parallelizable
- Very good for our problems of computing norms of \mathbf{F}
- Can be readily extended to variations of the basic stream model

$$\text{polylog}(n, m) \begin{bmatrix} \text{sketch matrix} \end{bmatrix} = \begin{bmatrix} \vdots \end{bmatrix} \text{sketch vector}$$

\mathbf{F}

The diagram illustrates the computation of a linear sketch. On the left, a sketch matrix is shown as a red bracket containing the text "sketch matrix". Above this matrix is a red double-headed vertical arrow labeled "polylog(n, m)". To the right of the sketch matrix is an equals sign (=). To the right of the equals sign is a blue bracket containing a vertical vector with several horizontal segments and a middle dot, labeled "sketch vector". Above this vector is a black bracket containing the letter "F".

Data Stream Model: Time Series Model

Time Series Model

Every stream item gives the current frequency of an element ($\mathbf{F}[a_i]$)

Stream items are $a_i = \langle j, c_i \rangle$ and it means $\mathbf{F}[j] \leftarrow c_i$

For stream $\mathcal{S} : \langle 7, 3 \rangle, \langle 3, 3 \rangle, \langle 2, 9 \rangle, \langle 7, 2 \rangle, \langle 9, 1 \rangle, \langle 3, 1 \rangle$

The final frequency vector will be

$$\mathbf{F} = \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \hline 0 & 9 & 1 & 0 & 0 & 0 & 2 & 0 & 1 \\ \hline \end{array}$$

- Used to measure link-bandwidth or energy consumption over time
- Very useful if there are multiple streams (e.g. stock prices for different companies)

Data Stream Model: Cash-Register Model

Cash-Register Model aka Arrivals-Only Stream

Every stream item is an increment to a frequency.

Stream items are $a_i = \langle j, c_i \rangle$ and it means $\mathbf{F}[j] \leftarrow \mathbf{F}[j] + c_i$ $c_i \geq 1$

For stream $S : \langle 7, 3 \rangle, \langle 3, 3 \rangle, \langle 2, 9 \rangle, \langle 7, 2 \rangle, \langle 9, 1 \rangle, \langle 3, 1 \rangle$

The final frequency vector will be

$$\mathbf{F} = \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \hline 0 & 9 & 4 & 0 & 0 & 0 & 5 & 0 & 1 \\ \hline \end{array}$$

Can be used e.g. for packet counts in every flow

Data Stream Model: Turnstile Model

Turnstile Model aka Arrivals and Departures Stream

Every stream item is an update to a frequency

Stream items are $a_i = \langle j, c_i \rangle$ and it means $\mathbf{F}[j] \leftarrow \mathbf{F}[j] + c_i$ $c_i \geq 1$

For stream $S : \langle 7, 3 \rangle, \langle 3, 3 \rangle, \langle 2, 9 \rangle, \langle 7, -2 \rangle, \langle 9, 1 \rangle, \langle 3, -1 \rangle$

The final frequency vector will be

$$\mathbf{F} = \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \hline 0 & 9 & 2 & 0 & 0 & 0 & 1 & 0 & 1 \\ \hline \end{array}$$

Generally, model has restriction of $\mathbf{F}[\cdot] \geq 0$

Universal hash functions

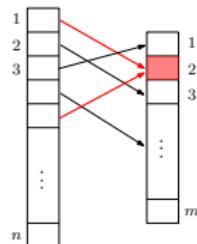
Hash functions/table is an efficient way to implement the Dictionary ADT

Hash functions map keys $A \subset U$ to m buckets labeled $\{0, 1, 2, \dots, m - 1\}$

A is not known in advance and $|A| = n$

Desired properties from hashing

- Fewer collisions
- Small range (m)
- Small space complexity to store hash function
- Easy to evaluate hash value for any key



A family of hash functions \mathcal{H} is 2-universal if

$$\text{for any distinct keys } x, y \in U, \quad \Pr_{h \in R \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{m}$$

Source of randomness is picking h (at random) from the family

Universal hash functions

A family of hash functions \mathcal{H} is 2-universal if

$$\text{for any distinct keys } x, y \in U, \quad \Pr_{h \in_R \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{m}$$

Linear Congruential Generators for $U = \mathbb{Z}$

- Pick a prime number $p > m$
- For any two integers a and b ($1 \leq a \leq p - 1$), ($0 \leq b \leq p - 1$)
- A hash function $h_{a,b} : U \mapsto [m]$ is defined as

$$h_{a,b}(x) = (ax + b) \pmod{p} \pmod{m}$$

$\mathcal{H} := \{h_{a,b} : 1 \leq a \leq p - 1, 0 \leq b \leq p - 1\}$ is 2-universal

Picking a random $h \in \mathcal{H}$ amounts to picking random a and b

Count-Min Sketch

- Count-Min sketch (Cormode & Muthukrishnan 2005) for frequency estimates
- Cannot store frequency of every elements
- Store total frequency of random groups (elements in hash buckets)

Algorithm : Count-Min Sketch (k, ϵ, δ)

COUNT \leftarrow ZEROS(k) ▷ sketch consists of k integers

Pick a random $h : [n] \mapsto [k]$ from a 2-universal family \mathcal{H}

On input a_i

COUNT[$h(a_i)$] \leftarrow COUNT[$h(a_i)$] + 1 ▷ increment count at index $h(a_i)$

On query j

▷ query: $F[j] = ?$

return COUNT[$h(j)$]

Count-Min Sketch

Algorithm : Count-Min Sketch (k, ϵ, δ)

COUNT \leftarrow ZEROS(k) ▷ sketch consists of k integers

Pick a random $h : [n] \mapsto [k]$ from a 2-universal family \mathcal{H}

On input a_i

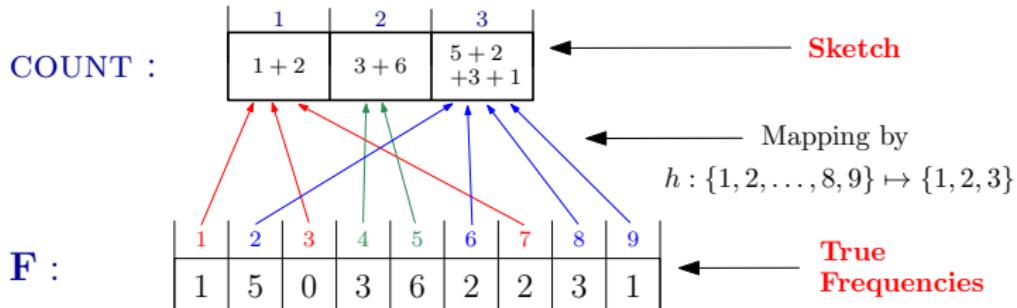
COUNT[$h(a_i)$] \leftarrow COUNT[$h(a_i)$] + 1 ▷ increment count at index $h(a_i)$

On query j

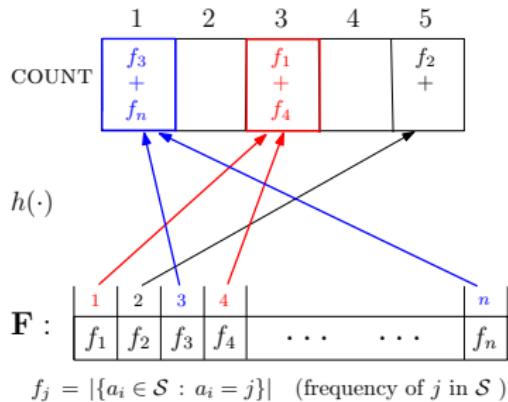
▷ query: $F[j] = ?$

return COUNT[$h(j)$]

$$\mathcal{S} : 2, 5, 6, 7, 8, 2, 1, 2, 7, 5, 5, 4, 2, 8, 8, 9, 5, 6, 4, 4, 2, 5, 5$$



Count-Min Sketch

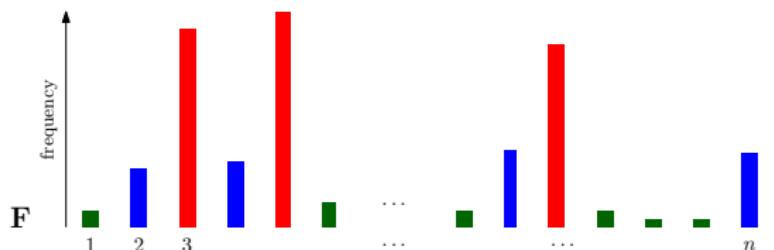


- $k = \frac{2}{\epsilon}$
- Large k means better estimate (smaller groups) but more space
- \tilde{f}_j : estimate for f_j – output of algorithm

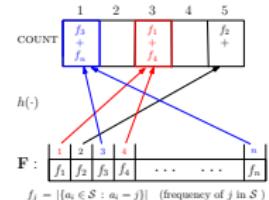
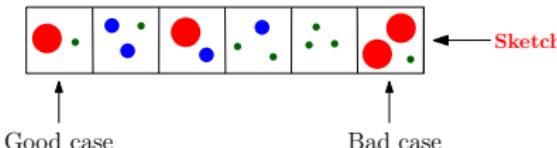
Count-Min Sketch

- $k = 2/\epsilon$
- Large k means better estimate but more space
- \tilde{f}_j : estimate for f_j – output of algorithm

Bounds on \tilde{f}_j : (idea)



COUNT :



$h(\cdot)$

\mathbf{F} :

$f_j = |\{a_i \in \mathcal{S} : a_i = j\}|$ (frequency of j in \mathcal{S})

Count-Min Sketch

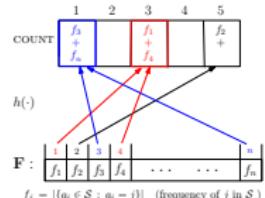
- $k = 2/\epsilon$
- Large k means better estimate but more space
- \tilde{f}_j : estimate for f_j – output of algorithm

Bounds on \tilde{f}_j : (idea)

- 1 $\tilde{f}_j \geq f_j$
 - Other elements that hash to $h(j)$ contribute to \tilde{f}_j
- 2 $Pr[\tilde{f}_j \leq f_j + \epsilon \|F\|_1] \geq \frac{1}{2}$
 - $X_j = \tilde{f}_j - f_j$ ▷ Excess in \tilde{f}_j (error)
 - $X_j = \sum_{i \in [n] \setminus j} f_i \cdot 1_{h(i)=h(j)}$ ▷ $1_{\text{condition}}$ is indicator of condition

$$\mathbb{E}(X_j) = \mathbb{E}\left(\sum_{i \in [n] \setminus j} f_i \cdot 1_{h(i)=h(j)}\right) = \sum_{i \in [n] \setminus j} f_i \cdot \frac{1}{k} \leq \sum_{i \in [n] \setminus j} \|F\|_1 \cdot \frac{\epsilon}{2}$$

- By Markov inequality we get the bound



Count-Min Sketch

Idea: Amplify the probability of the basic count-min sketch

Keep t over-estimates, $t = \log(1/\delta)$, $k = 2/\epsilon$ and return their minimum

Unlikely that all t functions hash j with very frequent elements

Algorithm : Count-Min Sketch (k, ϵ, δ)

COUNT \leftarrow ZEROS($t \times k$) ▷ sketch consists of t rows of k integers

Pick t random functions $h_1, \dots, h_t : [n] \mapsto [k]$ from a 2-universal family

On input a_i :

for $r = 1$ to t **do**

COUNT[r][$h_r(a_i)$] \leftarrow COUNT[r][$h_r(a_i)$] + 1

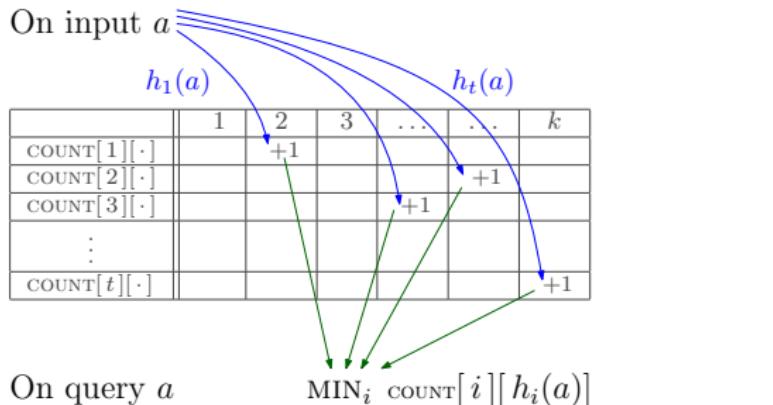
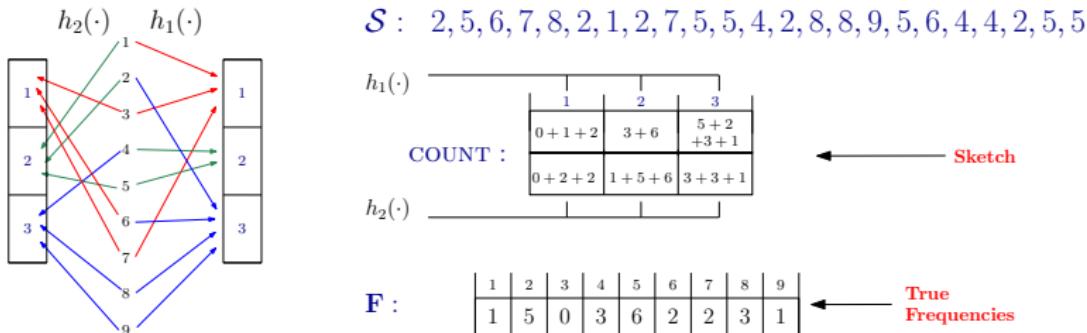
▷ increment COUNT[r] at index $h_r(a_i)$

On query j

▷ query: F[j] = ?

return $\min_{1 \leq r \leq t} \text{COUNT}[r][h_r(j)]$

Count-Min Sketch



Count-Min Sketch

1 $\tilde{f}_j \geq f_j$

- For every r , other elements that hash to $h_r(j)$ contribute to \tilde{f}_j

2 $\tilde{f}_j \leq f_j + \epsilon \|F\|_1$ with probability at least $1 - \delta$

- X_{jr} : contribution of other elements to $Count[r][h_r(j)]$
- $\Pr[X_{jr} \geq \epsilon \|F\|_1] \leq \frac{1}{2}$ for $k = 2/\epsilon$
- The event $\tilde{f}_j \geq f_j + \epsilon \|F\|_1$ is $\forall 1 \leq r \leq t \quad X_{jr} \geq \epsilon \|F\|_1$
- $\Pr[\forall r X_{jr} \geq \epsilon \|F\|_1] \leq \left(\frac{1}{2}\right)^t$
- $t = \log(\frac{1}{\delta}) \implies \Pr[\forall r X_{jr} \geq \epsilon \|F\|_1] \leq \left(\frac{1}{2}\right)^{\log^{1/\delta}} = \delta$

- Count-Min sketch is an $(\epsilon \|F\|_1, \delta)$ -additive approximation algorithm
- Space required is $k \cdot t$ integers = $O(1/\epsilon \log(1/\delta) \log n)$ (plus constant)

The Count Sketch

- In Count-Min sketch error in frequency estimate accumulates (group total)
- The Count Sketch (Charikar, Chen, Farach-Colton, 2002)
- A frequency estimate where errors in a group cancel each other

Algorithm : Count Sketch (k, ϵ, δ)

Pick a random $h : [n] \mapsto [k]$ from a 2-universal family \mathcal{H}

Pick a random $g : [n] \mapsto \{-1, 1\}$ from a 2-universal family

$\text{COUNT} \leftarrow \text{ZEROS}(k)$ ▷ sketch consists of k integers

On input a_i

$\text{COUNT}[h(a_i)] \leftarrow \text{COUNT}[h(a_i)] + g(a_i)$

▷ increment or decrement, depending on value of $g(a_i)$ COUNT at index $h(a_i)$

On query j

▷ query: $\mathbf{F}[j] = ?$

return $g(j) \times \text{COUNT}[h(j)]$

Count Sketch

Algorithm : Count Sketch (k, ϵ, δ)

Pick a random $h : [n] \mapsto [k]$ from a 2-universal family \mathcal{H}

Pick a random $g : [n] \mapsto \{-1, 1\}$ from a 2-universal family

COUNT \leftarrow ZEROS(k)

▷ sketch consists of k integers

On input a_i :

COUNT[$h(a_i)$] \leftarrow COUNT[$h(a_i)$] + $g(a_i)$

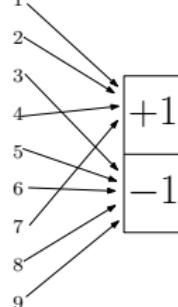
▷ increment or decrement, depending on value of $g(a_i)$ COUNT at index $h(a_i)$

On query j

▷ query: $F[j] = ?$

return $g(j) \times \text{COUNT}[h(j)]$

$g(\cdot)$



$S : 2, 5, 6, 7, 8, 2, 1, 2, 7, 5, 5, 4, 2, 8, 8, 9, 5, 6, 4, 4, 2, 5, 5$

COUNT :

1	2	3
$+1 + 2$	$+3 - 6$	$+5 - 2$ $-3 - 1$

Sketch

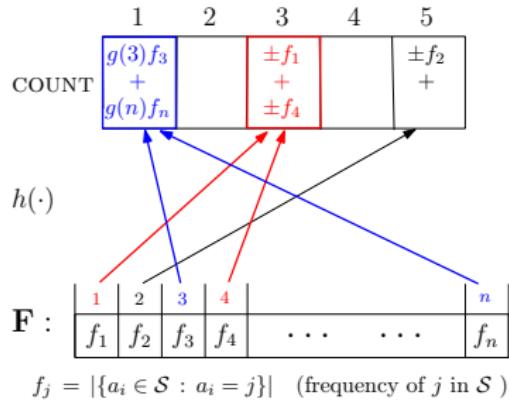
$F :$

1	5	0	3	6	2	2	3	1
---	---	---	---	---	---	---	---	---

Mapping by
 $h : \{1, 2, \dots, 8, 9\} \mapsto \{1, 2, 3\}$

True Frequencies

The Count Sketch



- $k = \frac{3}{\epsilon^2}$
- \tilde{f}_j : estimate for f_j – output of algorithm

The Count Sketch

- $k = 3/\epsilon^2$
- \tilde{f}_j : estimate for f_j – output of algorithm

Bounds on \tilde{f}_j :

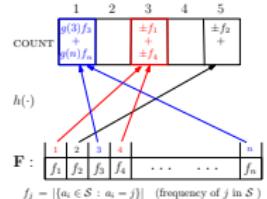
$$1 \quad E(\tilde{f}_j) = f_j$$

$$\text{COUNT}[h(j)] = \sum_{i \in [n]} f_i \cdot g(i) \cdot 1_{h(i)=h(j)}$$

$$\tilde{f}_j = g(j) \sum_{i \in [n]} f_i \cdot g(i) 1_{h(i)=h(j)} = g(j) \left(f(j)g(j) + \sum_{i \in [n] \setminus j} f_i \cdot g(i) 1_{h(i)=h(j)} \right)$$

$$= f(j)(g(j))^2 + \sum_{i \in [n] \setminus j} f_i \cdot g(i)g(j) \cdot 1_{h(i)=h(j)} = f(j) + \sum_{i \in [n] \setminus j} f_i \cdot g(i)g(j) 1_{h(i)=h(j)}$$

$$\implies \mathbb{E}(\tilde{f}_j) = f_j \quad \triangleright \mathbb{E}(1_{h(i)=h(j)}) = \frac{1}{k} \text{ and } \mathbb{E}(g(i)g(j)) = 0$$



The Count Sketch

- $k = 3/\epsilon^2$
- \tilde{f}_j : estimate for f_j – output of algorithm

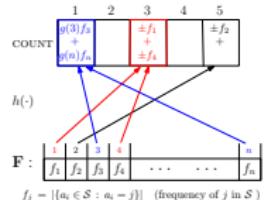
Bounds on \tilde{f}_j :

1 $E(\tilde{f}_j) = f_j$

2 $Var(\tilde{f}_j) \leq \frac{1}{k} \|F\|_2$

3 $Pr[|\tilde{f}_j - f_j| \geq \epsilon \|F\|_2] \leq 1/3$

- substitute $k = 3/\epsilon^2$ and use Chebychev inequality



Count Sketch

Probability Amplification

Algorithm : Count Sketch (k, ϵ, δ)

COUNT \leftarrow ZEROS($t \times k$) ▷ sketch consists of t rows of k integers

Pick t random functions $h_1, \dots, h_t : [n] \mapsto [k]$ from a 2-universal family

Pick t random functions $g_1, \dots, g_t : [n] \mapsto \{-1, 1\}$ from a 2-uni. family

On input a_i

for $r = 1$ to t **do**

COUNT[r][$h_r(a_i)$] \leftarrow COUNT[r][$h_r(a_i)$] + $g_r(a_i)$

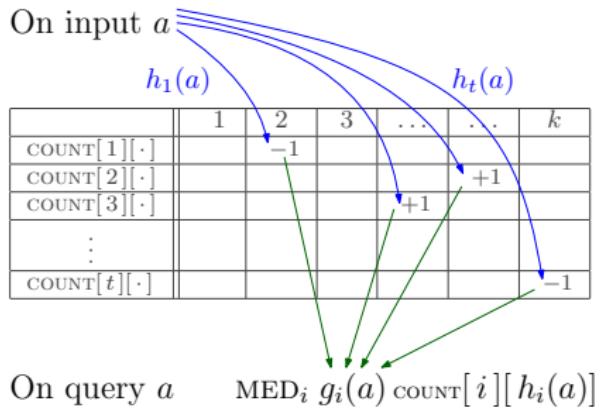
▷ inc/dec COUNT[r] at index $h_r(a_i)$

On query j ▷ query: F[j] = ?

return MEDIAN _{$1 \leq r \leq t$} $g_r(j) \times$ COUNT[r][$h_r(j)$]

Count Sketch

Keep t unbiased estimates, $t = \log(1/\delta)$, $k = 3/\epsilon^2$. Their median is a good estimate, unless at least $t/2$ estimates are very bad



- 1 $E(\tilde{f}_j) = f_j$
- 2 $|\tilde{f}_j - f_j| \leq \epsilon \|F\|_2$ with probability at least $1 - \delta$ ▷ **Uses Chernoff bound**
- Count sketch is an $(\epsilon \|F\|_2, \delta)$ additive approximation algorithm
- Space required is $k \cdot t$ integers = $O(1/\epsilon^2 \log(1/\delta) \log n)$ (plus constant)

Estimate F_2 : AMS Algorithm

- 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} $\mathbf{F} = \{f_1, f_2, \dots, f_n\}$

$$F_2 := \sum_{i=1}^n f_i^2 \quad \triangleright \text{second frequency moment}$$

Easy to compute if we store F $\triangleright O(n)$ space

Can store $f_1 + f_2 + \dots + f_n$ $\triangleright O(1)$ space

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

Estimate F_2 : AMS Algorithm

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

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

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

Algorithm:

for each $a_i \in S$

$X \leftarrow X + 1$

return X^2

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

Estimate F_2 : AMS Algorithm

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

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

$$(f_1 + f_2 + f_3 + f_4)^2 = \underline{f_1^2 + f_2^2 + f_3^2 + f_4^2} + 2(f_1 f_2 + f_1 f_3 + f_2 f_3 + f_1 f_4 + f_2 f_4 + f_3 f_4)$$

$$(f_1 - f_2 + f_3 - f_4)^2 = \underline{f_1^2 + f_2^2 + f_3^2 + f_4^2} + 2(-f_1 f_2 + f_1 f_3 - f_2 f_3 - f_1 f_4 + f_2 f_4 - f_3 f_4)$$

Estimate F_2 : AMS Algorithm

Algorithm (AMS):

$g : [n] \rightarrow \{-1, +1\}$

▷ random hash function

for each $a_i \in \mathcal{S}$

$X \leftarrow X + g(a_i)$

return X^2

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

Estimate F_2 : AMS Algorithm

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

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

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

Estimate F_2 : AMS Algorithm

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

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

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

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

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

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

Estimate F_2 : AMS Algorithm

Algorithm (AMS):

$g : [n] \mapsto \{-1, +1\}$

▷ random hash function

for each $a_i \in S$

$$X \leftarrow X + g(a_i)$$

return X^2

$$\mathbb{E}[X^2] = F_2 \quad \text{Var}(X^2) \leq 2F_2$$

Amplify the probability

- Keep $k = 2/\epsilon^2\delta$ counters
- Return their average $\bar{X} = \frac{1}{k} \sum_i X_i^2$

$$\Pr[|\bar{X} - F_2| > \epsilon F_2] \leq \delta$$

Estimate F_2 : AMS Algorithm

Algorithm (AMS):

$g : [n] \mapsto \{-1, +1\}$

▷ random hash function

for each $a_i \in \mathcal{S}$

$$X \leftarrow X + g(a_i)$$

return X^2

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

$$\mathbf{g} = \boxed{g(1) \quad g(2) \quad \dots \quad g(n)}$$

$$\begin{matrix} \mathbf{F} \\ \boxed{f_1} \\ \boxed{f_2} \\ \vdots \\ \vdots \\ \boxed{f_n} \end{matrix} = X$$

Estimate F_2 : AMS Algorithm

Algorithm (AMS):

$$g : [n] \mapsto \{-1, +1\}$$

▷ random hash function

for each $a_i \in \mathcal{S}$

$$X \leftarrow X + g(a_i)$$

return X^2

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

$$\mathbf{g} = \boxed{+1 \quad -1 \quad \dots \quad \quad +1}$$

$$\begin{matrix} \mathbf{F} \\ \boxed{f_1} \\ \boxed{f_2} \\ \vdots \\ \vdots \\ \boxed{f_n} \end{matrix} = X$$

Estimate F_2 : AMS Algorithm

Algorithm (AMS):

$g : [n] \mapsto \{-1, +1\}$

▷ random hash function

for each $a_i \in \mathcal{S}$

$$X \leftarrow X + g(a_i)$$

return X^2

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

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

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

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

Estimate F_2 : AMS Algorithm

Algorithm (AMS):

$g : [n] \mapsto \{-1, +1\}$

▷ random hash function

for each $a_i \in \mathcal{S}$

$$X \leftarrow X + g(a_i)$$

return X^2

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

$\mathbf{G} =$

+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

Estimate F_2 : AMS Algorithm

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

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

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

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

With probability at least $1 - \delta$

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

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

Estimate F_2 : AMS Algorithm

$$\mathbf{G} = \begin{bmatrix} +1 & -1 & \dots & & +1 \\ -1 & -1 & \dots & & -1 \\ \vdots & & & & \vdots \\ -1 & +1 & \dots & & -1 \end{bmatrix}$$

$$\mathbf{F} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_k \end{bmatrix}$$

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

\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 $U, V \in \mathcal{R}^n$

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

Johnson-Lindenstrauss Lemma

- Given $V = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} \subset \mathcal{R}^d$
- For any $\epsilon \in (0, 1/2)$, there is a linear map $f : \mathcal{R}^d \mapsto \mathcal{R}^k$
- $k = c \log n / \epsilon^2$, such that for any $\mathbf{u}, \mathbf{v} \in V$

$$(1 - \epsilon) \|\mathbf{u} - \mathbf{v}\|_2 \leq \|f(\mathbf{u}) - f(\mathbf{v})\|_2 \leq (1 + \epsilon) \|\mathbf{u} - \mathbf{v}\|_2$$

- This map can be obtained very easily
- Let \mathbf{M} be a $k \times d$ matrix, with $M_{ij} \in \mathcal{N}(0, 1)$, then

$$f(\mathbf{u}) = \frac{1}{\sqrt{k}} \mathbf{M} \mathbf{u}$$