

# Contents

<b>1</b>	<b>The data stream phenomenon</b>	<b>1</b>
1.1	Characteristic of data streams . . . . .	2
1.2	Data Stream Processing Model . . . . .	3
1.2.1	Approximation Algorithms . . . . .	3
<b>2</b>	<b>Data Stream Applications</b>	<b>4</b>
2.1	Sensor Data . . . . .	5
2.2	Network Monitoring and Management . . . . .	5
2.3	Web Click Stream Analysis . . . . .	6
2.4	Search Queries Stream . . . . .	6
2.5	Financial Time Series . . . . .	7
2.6	Optimized Query Execution Plan . . . . .	7
<b>3</b>	<b>Stream Model of Computation</b>	<b>8</b>
<b>4</b>	<b>Synopsis of the stream</b>	<b>9</b>
4.0.1	Synopsis Based Exact Stream Computation . . . . .	9
4.1	Finding Missing Integer(s): A motivating example . . . . .	10
4.2	Synopsis: Sliding Window . . . . .	10
4.3	Synopsis: Random Sample . . . . .	11
4.4	Synopsis: Histogram . . . . .	13
<b>5</b>	<b>Models of data streams</b>	<b>13</b>
5.1	Synopsis: Linear Sketch . . . . .	14
<b>6</b>	<b>Finding frequency of an element: The Count-Min sketch</b>	<b>15</b>
6.1	Applications . . . . .	16
6.2	The Count-Min sketch . . . . .	16
6.2.1	Amplifying the probability . . . . .	19
<b>7</b>	<b>Finding frequency of an element: The Count sketch</b>	<b>20</b>
7.1	The Count sketch . . . . .	20
7.1.1	Amplifying the probability . . . . .	24
<b>8</b>	<b>Comparison of the two sketches</b>	<b>26</b>
<b>9</b>	<b>Lower bound</b>	<b>26</b>
9.1	Index problem . . . . .	27
9.2	Reduction from the Frequency estimation to Index problem . . . . .	27

## 1 The data stream phenomenon

**Stream Processing or Stream Analytics or Streaming Algorithms** is the application of data analysis and algorithmic methods on a continuous stream of data items and drawing meaningful analytics and knowledge from it.

Stream processing is generally performed in the following settings.

- Single Pass over the stream: The common requirement is to process each item exactly once. Though in certain cases we do use 2 pass or a small number of passes over the stream.
- Limited Memory: The common characteristic of stream processing is the requirement to use poly-logarithmic space (in length of stream or domain of data objects). This immediately implies that we cannot store the whole stream (and the single pass requirement becomes clear)
- Constant per item processing: Given an object the requirement is to process it in near real time.
- Arbitrary arrival order: Generally, in streaming algorithm we do not make any assumption on distribution or order of items. For instance we cannot say that objects are streamed in a random order or that objects are coming from a certain distribution. In other words, we essentially assume adversarial (worse-case) order.

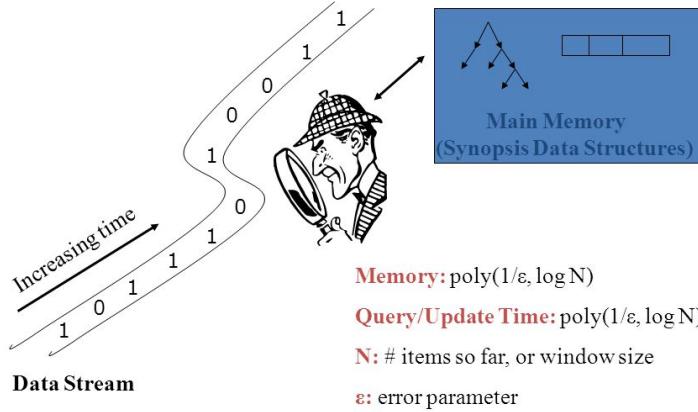
## 1.1 Characteristic of data streams

Some characteristic of data streams or some key aspects in which data stream processing is different than the usual data processing are given as follows. These are based on the textbook Han & Kamber, Data Mining Concepts & Techniques, 2nd Ed.

- 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

The following is a contrast with traditional data analysis and in what aspects stream data is fundamentally different than traditional datasets based on a tutorial by Rajeev Motwani in PODS 2001.

Motwani, PODS (2002)



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

## 1.2 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)

### 1.2.1 Approximation Algorithms

As mentioned above, in most of interesting cases of computing functions over streams,  $f(\mathcal{S})$  can provably not be computed given the sublinear space and 1 pass requirements. We, therefore often allow **approximation algorithms** that make errors with bounded probability.

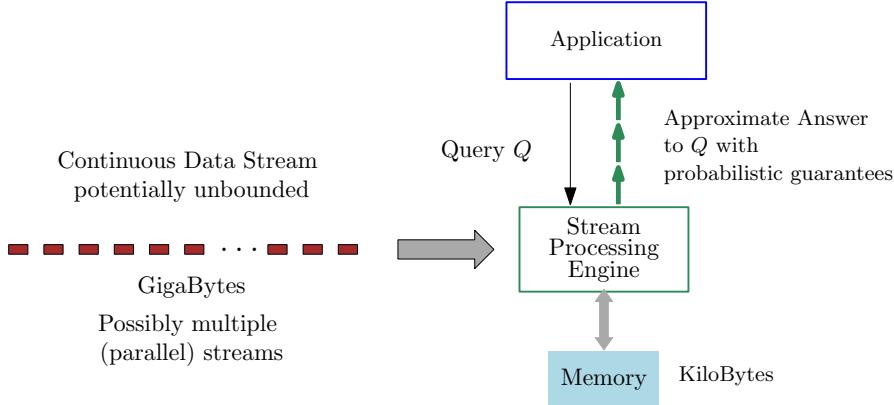
**Definition 1.** Let  $\mathcal{A}$  be an algorithm to compute  $f(\mathcal{S})$ . Denote by  $\mathcal{A}(\mathcal{S})$  the output of  $\mathcal{A}$  on input  $\mathcal{S}$ . We say that  $\mathcal{A}$  is an  $(\epsilon, \delta)$  approximation algorithm if

$$\mathbb{P}((\mathcal{A}(\mathcal{S}) - f(\mathcal{S})) > \epsilon f(\mathcal{S})) \leq \delta.$$

Instead of the above requirement of a multiplicative approximation guarantee, an algorithm could provide an additive approximation guarantee.

**Definition 2.** Let  $\mathcal{A}$  be an algorithm to compute  $f(\mathcal{S})$ . Denote by  $\mathcal{A}(\mathcal{S})$  the output of  $\mathcal{A}$  on input  $\mathcal{S}$ . We say that  $\mathcal{A}$  is an  $(\epsilon, \delta)$  additive approximation algorithm if

$$\mathbb{P}((\mathcal{A}(\mathcal{S}) - f(\mathcal{S})) > \epsilon) \leq \delta.$$

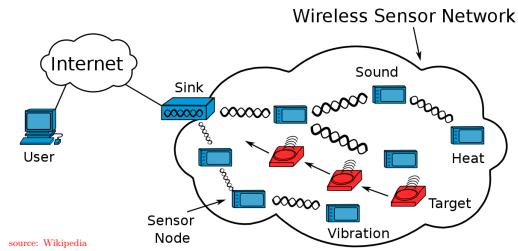


## 2 Data Stream Applications

Stream data comes in many domains and has various applications. Here are some application domains from the textbook Han & Kamber, Data Mining Concepts & Techniques, 2nd Ed..

- 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)

We elaborate on some of these applications.



## 2.1 Sensor Data

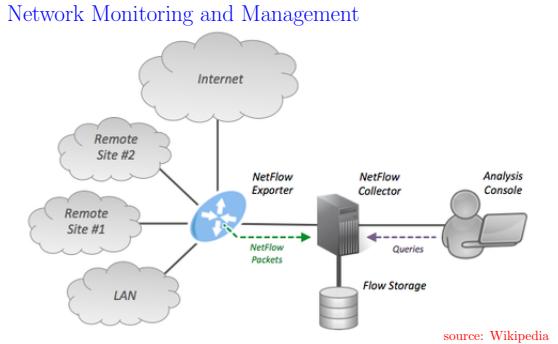
With the advent of IoT and sensor networks, there has been a huge amount of data generated by sensor nodes. The data is practically unlimited but the sensor devices have very limited computation power and memory. The limited computation power and (battery) power requirements to transmit data necessitate processing the stream of data at the sensor node without the need to storing or transmitting the whole data. It is estimated that 1 bit transmission consumes power almost equivalent to executing 800 instructions (Madden et.al. (2002)). Certain key statistical aggregates and drawn analytics can be transmitted to a server for detailed processing. We note that even if there is no storage or communication bandwidth limitation, the near real time requirements for analytics warrants stream processing, while can keep the whole dataset (in a persistent manner) at the back-end or further detailed processing. Results of the basic stream analytics can be used to guide detailed and sophisticated analyses offline in data warehouses.

## 2.2 Network Monitoring and Management

**NetFlow:** A Cisco tool for network administrators that is widely used for performance metrics, security analysis, detection and forensics.

For each Flow it reports (logs)

- Network Interface
- Source/Destination IP Addresses
- IP Protocol
- Source/Destination port
- TCP Flags
- Total packets/bytes in flow,
- and much more



Stream of this NetFlow data is widely used for network traffic engineering and monitoring, Traffic Volume estimation & analysis, Load balancing and Efficient Resource Utilization. This is also used for various network security purposes such as (D)DOS Attack Detection and identifying Service Level Agreement( SLA) Violation in real time [4]

Following are some example queries that can be (approximately) answered in the streaming setting (single pass and sub-linear memory). It is not very hard to see how answers to these queries can and will be used for the applications mentioned above.

- How many bytes sent b/w IP-1 and IP-2?

- How many IP addresses are active?
- What are the top 100 IP's by traffic volume (busiest)?
- What is the average duration of IP session and standard deviation?
- What is the median number of bytes in an IP session?
- Find sessions that transmitted more than  $x$  bytes
- Find sessions whose duration is more than twice the average session duration,
- List all IP's with a sudden spike in traffic
- List all IP involved in more than 1k sessions

As an example, AT&T Processes over 567 billion flow records per day, which amounts to about  $\sim 15$  PBytes of data (Fred Stinger (AT&T) FloCon (2017) Netflow Collection and Analysis ..). AT&T detects and characterizes approximately 500 anomalies per day using this data.

## 2.3 Web Click Stream Analysis



Figure credit: Alex Smola @Yahoo research & ANU

- Stream of user clicks on websites (tracked via cookies)
- Find hot links, frequent IP's, click probability
- Enhanced customer experience & conversion rates
- Digital marketing – Up-selling and cross-selling



## 2.4 Search Queries Stream

- Discover trends and patterns
- Relevant keywords for website
- Estimate competition scores or difficulty
- Estimate keywords CPC (cost per click)





### AMI Data Streams

Energy consumption Analysis from smart meters transmitting information of load per minute [3, 1, 2].



- Electricity consumption data from AMI (Automatic Metering Interface)
- Find average hourly load, load surges, anomaly
- Short term load forecast (total or for individual consumer)
- Identify faults, drops, failures

## 2.5 Financial Time Series

- Time stamped real time (multiple) stock data
- Need near real time prediction
- Algorithmic Trading

## 2.6 Optimized Query Execution Plan

In some cases the data doesn't have to be streamed in, but one can deal with massive data that is stored on disk in a streaming fashion (Single Pass, Sequential I/O). Such algorithms are also judged by the per item processing time, storage required and computation time.



Suppose we have a large students database with many fields. We do not have the usual (sub-linear) space restriction, streaming could still help us. We give an example of optimizing query execution plans by *selectivity estimation*. Suppose we want to execute a query like

```
SELECT * FROM Table WHERE 25 <= age <= 35 AND 54 <= weight <= 60
```

Suppose there are  $n = 1M$  students in total. **Brute force execution time** is  $2n$  comparisons, as we have to compare the two fields of every record.

If we have some statistical information about the distribution of each attribute in particular age and weight, then runtime will be significantly different.

Age	Freq	Weight	Freq.
0 – 10	7%	0 – 20	20%
11 – 20	8%	21 – 40	25%
21 – 30	10%	41 – 60	10%
31 – 40	12%	61 – 80	15%
41 – 50	13%	81+	30%
51 – 60	25%		
61 – 70	20%		
71+	5%		

If we first filter on Age, then on weight i.e. first compare every records age to see if it is within the range and then compare the weight of only the filtered records, then the query can be executed in using  $1.22n$  comparisons. On the other hand first filtering on weight, then on age would require  $1.1n$  comparisons.

These aggregate statistics can be obtained using one scan on the whole database for all attributes or while the data is being stored in the database. These kind of techniques are often used for executing joins though (not just select queries). These techniques can also be used for performing Exploratory data analytics or Database monitoring.

### 3 Stream Model of Computation

A data stream is a sequence of  $m$  items  $\mathcal{S} = a_1, a_2, a_3, \dots, a_m$ , where each item is chosen from some universe of size  $n$ . Typically we take the universe to be  $[n] := \{1, 2, \dots, n\}$ . Some comments

about streams:

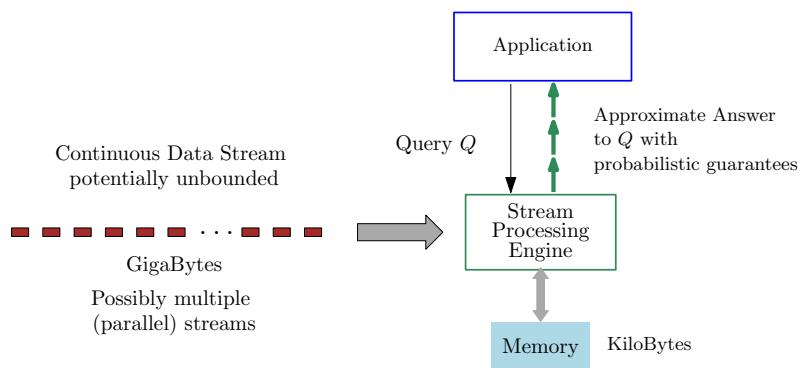
- $n$  and  $m$  are two size parameters.
- We do not assume anything about **distribution of items**
- Typically we work in the model where we see each **item only once** in the order given by  $\mathcal{S}$ , in general we cannot (or do not want to) save the whole stream. In the literature there are algorithms that take more than one passes over the stream but here we restrict ourselves to one-pass algorithms only.
- Our goal is to use a very **small amount of memory** for computing some function over stream,  $f(\mathcal{S})$ , i.e. space requirement of algorithm should be  $o(\min\{m, n\})$ . Ideally we should use  $O(\log n + \log m)$  space. Note that this space is needed anyway to store one (or some constant number of) stream items for processing. Sometimes space requirement could be relaxed to polylogarithmic in  $\min\{m, n\}$  (like  $(\log n)^c$  for some constant  $c$ ).
- In most of interesting cases of computing functions over streams,  $f(\mathcal{S})$  can provably not be computed given the sublinear space and 1 pass requirements. We, therefore often allow **approximation algorithms** that make errors with bounded probability.

## 4 Synopsis of the stream

The Fundamental Methodology in stream processing is to keep a synopsis of the stream and answer query based on it. The synopsis is updated after examining each item in  $O(1)$  time. The synopsis is a succinct summary of the stream (observed so far) with total size poly-log bits

Following are some major families of synopsis.

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



### 4.0.1 Synopsis Based Exact Stream Computation

Given a stream  $S = a_1, a_2, \dots, a_m$ , such that each  $a_i \in [n]$ , we can compute

- Length of  $\mathcal{S}$  ( $m$ ) by storing a running counter
- Sum of  $\mathcal{S}$  by storing a running sum
- Mean of  $\mathcal{S}$  from sum and length of  $\mathcal{S}$

- Variance of  $\mathcal{S}$  from sum, sum of square, and length of  $\mathcal{S}$

$$Var(X) = E(X^2) - (E(X))^2$$

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

This is just to emphasize the fact, that though the streaming model is restrictive, we can still achieve quite a lot.

## 4.1 Finding Missing Integer(s): A motivating example

Suppose you are given  $n - 1$  distinct positive integer in  $[n]$  and you want to find the one integer that is missing

You can solve this problem easily by maintaining a bit vector  $B[1 \dots n]$  and marking  $B[i] = 1$  if the input integer is  $i$ . Then at the end of input, one scan through  $B$  and identifying the index of 0 reveals the missing integer.

Runtime of this algorithm is  $2n - 1$  and space is  $n$  bits.

You can accomplish this by realizing that  $\sum_{i=1}^n i = \frac{n(n+1)}{2}$  and maintaining only one integer  $S$  that is the running sum of all input integers.

Then the missing integer is  $\frac{n(n+1)}{2} - S$

Runtime of this algorithm is  $n$  additions and space complexity is  $2n \log n$

What if two integers are missing and input is  $n - 2$  distinct integers in  $[n]$ . You can do this by maintaining the sum  $S$  and product  $P$  of input integers and in the end solving a system of two equations and two unknowns. i.e.

$$\frac{n(n+1)}{2} - S \quad \text{and} \quad n! - P$$

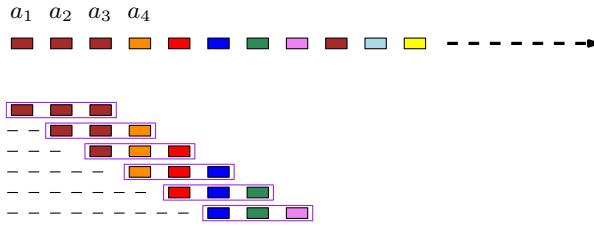
Though this will solve the problem, but observe that the product  $P$  of  $O(n)$  integers in  $[n]$  is  $n! \in \Omega(n^n)$  (by Sterling's approximation), hence storing  $P$  will take  $n \log n$  bits, which is about the same number of bits to store the whole input, (each integer taking  $\log n$  bit and there are  $n$  of them). Note that the original brute-force solution requires only  $n$  bits but that was two pass algorithm.

Another method, (**to find any constant  $k$  missing integers**) which is more space efficient is to maintain  $S_i$ , where  $S_i = \sum_{j \in \text{INPUT}} j^i$ , i.e. the sum of  $i^{th}$  powers of input integers. Since we know  $\sum_{i=1}^n j^i$ , in the end just as above we solve  $k$  equations with  $k$  unknowns to finding the  $k$  missing integers.

Notice that the above puzzle and the algorithms for mean, variance, length etc. are data stream algorithms that are exact and deterministic. Such algorithms are uncommon in the streaming model, generally randomized and approximate algorithms are known for problems in the streaming model. In the following we discuss some representative and widely applied problems.

## 4.2 Synopsis: Sliding Window

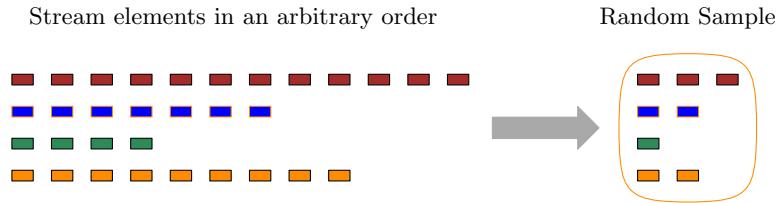
I am going to just mention the common families of synopses will complete them and add more material at some later time.



- 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

### 4.3 Synopsis: Random Sample

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

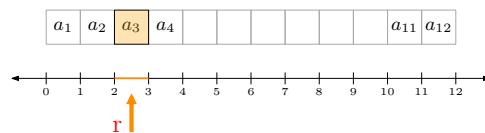


Sample a random element from array  $A$  of length  $n$

$\triangleright A[i]$  with prob  $1/n$

- Generate a random number  $r \in [0, n]$
- Return  $A[\lceil r \rceil]$

$\triangleright r \leftarrow \text{RAND}() \times n$

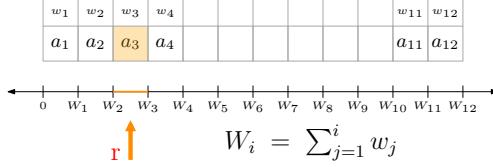


Sample random element (by weight) from array  $A$

$\triangleright A[i]$  with prob.  $w_i/W$

- Generate a random number  $r \in [0, \sum_{j=1}^n w_i]$
- Return  $A[i]$  if  $W_{i-1} \leq r < W_i$

$\triangleright r \leftarrow \text{RAND}() \times W_n$



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)

$$\begin{aligned}
&= \underbrace{\Pr \text{ that } a_i \text{ was selected at time } i}_{\frac{1}{i}} \times \underbrace{\Pr \text{ that } a_i \text{ survived in } R \text{ until time } m}_{\prod_{j=i+1}^m \left(1 - \frac{1}{j}\right)} \\
&= \cancel{\frac{1}{i}} \times \cancel{\frac{i}{i+1}} \times \cancel{\frac{i+1}{i+2}} \times \cancel{\frac{i+2}{i+3}} \times \dots \times \cancel{\frac{m-2}{m-1}} \times \cancel{\frac{m-1}{m}} = \frac{1}{m}
\end{aligned}$$

Sample  $k$  random elements from the stream  $S$  ▷  $a_i$  with prob.  $k/m$

---

**Algorithm** : Reservoir Sampling ( $\mathcal{S}, k$ )

---

```

 $R \leftarrow a_1, a_2, \dots, a_k$  ▷  $R$  (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)

$$\begin{aligned}
&= \underbrace{\Pr \text{ that } a_i \text{ was selected at time } i}_{\frac{k}{i}} \times \underbrace{\Pr \text{ that } a_i \text{ survived in } R \text{ until time } m}_{\prod_{j=i+1}^m \left(1 - \left(\frac{k}{j} \times \frac{1}{k}\right)\right)} \\
&= \cancel{\frac{1}{i}} \times \cancel{\frac{i}{i+1}} \times \cancel{\frac{i+1}{i+2}} \times \cancel{\frac{i+2}{i+3}} \times \dots \times \cancel{\frac{m-2}{m-1}} \times \cancel{\frac{m-1}{m}} = \frac{k}{m}
\end{aligned}$$

For more discussion on efficiently sampling  $k$  in streaming fashion with probability proportional to weights but with dynamic weights see [5].

## 4.4 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

## 5 Models of data streams

A data stream is a sequence of  $m$  items  $\mathcal{S} = a_1, a_2, a_3, \dots, a_m$ , where each item is chosen from some universe of size  $n$ . Typically we take the universe to be  $[n] := \{1, 2, \dots, n\}$ . We are interested in computing some statistical property of the multiset of items in  $\mathcal{S}$ . The stream describes some one-dimensional function (a signal), this is represented as an array of length  $n$ . Quite often the signal is the frequency vector  $F = [f_1 | f_2 | \dots | \dots | \dots | f_n]$  where  $f_i$  is the frequency of element  $i$ , i.e.  $f_j = |\{i : a_i = j\}|$ .

Various stream models differ how the items in the stream describe the signal. There are no standard definitions and a lot of variations. We consider the following three models and give examples of streams and the underlying signal.

1. **Time Series Model:** In this case if  $A$  is the underlying signal, then each  $a_i$  represents  $A[a_i]$ . For example if we want to record traffic on a given router every 5 minutes.  $A$  is an array of real number such that  $A[i]$  is the traffic at the  $i$ th 5-min internal and the the values of traffic reported by the router every 5 minutes are stream items ( $a_i$ ). Other example include the number of emails every day processed by a mail server, hourly trade volume in a certain stock exchange etc.

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

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

2. **Cash Register Model:** In this model each stream item is an increment to  $A[i]$ , so stream item  $a_i$  is like a pair  $(j, c)$ , ( $c > 0$ ) and it means that  $A_i[j] = A_{i-1}[j] + c$ . This is the most common model and is very applicable to the case of frequency, (so far we considered  $c = 1$  only). Applications: Suppose we want to monitor bandwidth consumed by different IP's. Each time an IP  $x$  uses a link it sends a certain number of packets, that could be treated as increments to  $x$ 's previous usage, CDR (call data record) is exactly this thing.

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

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

3. **Turnstile Model:** In this case  $a_i$ 's are updates to  $A[a_i]$  (instead of just increments we allow decrements also). Again think of  $a_i$  as a pair  $(j, c)$ , where  $c$  could be positive or negative. Usually it is very hard to get tight results in this model, that is one reason (certainly for us) not to study this. In some cases we restrict  $A[j]$  to always stay non-negative.

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

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

## 5.1 Synopsis: 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

$$\begin{array}{l} \mathcal{S} : a_1, a_2, a_3, a_4, \dots, a_m \\ a_i \in [n] \end{array} \quad \begin{array}{l} \mathbf{F} : \begin{array}{c|c|c|c|c|c|c|c|c|c} 1 & 2 & 3 & & & & & & n \\ \hline f_1 & f_2 & f_3 & & \dots & & \dots & & f_n \end{array} \\ f_j = |\{a_i \in \mathcal{S} : a_i = j\}| \text{ (frequency of } j \text{ in } \mathcal{S} \text{ )} \end{array}$$

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

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

$$\begin{array}{ll} \mathcal{S} = \langle a_1, a_2, a_3, \dots, a_m \rangle & a_i \in [n] \\ f_i : \text{frequency of } i \text{ in } \mathcal{S} & \mathbf{F} = \{f_1, f_2, \dots, f_n\} \end{array}$$

$$\begin{array}{ll} F_0 := \sum_{i=1}^n f_i^0 & \triangleright \text{number of distinct elements} \\ F_1 := \sum_{i=1}^n f_i & \triangleright \text{length of stream, } m \\ F_2 := \sum_{i=1}^n f_i^2 & \triangleright \text{second frequency moment} \end{array}$$

**Definition 3** (Linear Sketch). *Given a stream  $\mathcal{S}$ , a data structure  $sk(\mathcal{S})$  is called a sketch if*

- $sk(\mathcal{S})$  is computed in streaming fashion (by processing each item in the stream)
- We can easily combine the sketches of two streams  $\mathcal{S}_1$  and  $\mathcal{S}_2$ . More precisely, there is a space efficient streaming algorithm to compute  $sk(\mathcal{S}_1 \circ \mathcal{S}_2)$ , where  $\mathcal{S}_1 \circ \mathcal{S}_2$  is the concatenation of these two streams.

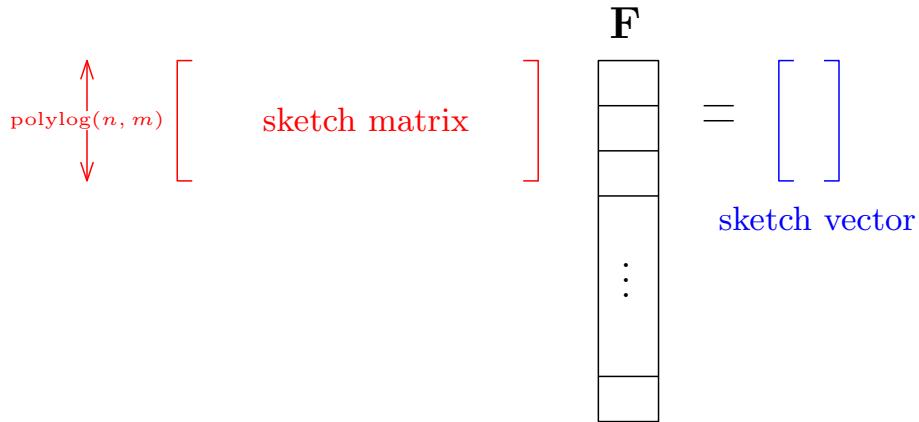
Basically (and quite often) we want to be able to efficiently update the sketch while processing the next item in the stream (think of  $\mathcal{S}_1 = a_1, a_2, \dots, a_{i-1}$  and  $\mathcal{S}_2 = a_i$ ). But technically since we can have multiple streams the definition takes care of it.

Linear sketch 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

## 6 Finding frequency of an element: The Count-Min sketch

we are given a data stream  $\mathcal{S} = a_1, a_2, a_3, \dots, a_m$ , where each  $a_i \in [n]$ . This implicitly defines a frequency vector  $F = \begin{array}{c|c|c|c|c|c|c|c|c|c} f_1 & f_2 & \dots & \dots & \dots & f_n \end{array}$  where  $f_i$  is the frequency of  $i$  in the stream, i.e.  $f_j = |\{i : a_i = j\}|$ . Our goal is to estimate the frequencies  $f_i$  for all elements.



## 6.1 Applications

One can easily imagine applications for such queries on a data stream in the cash-register model but with  $c$  fixed to be 1. It is an easy exercise to extend this to any increments values. We can trivially solve this by keeping a counter array of length  $n$ , but we want streaming solution.

For a quick application of estimates of frequency vectors we have used it to compute kernel values between two objects. Which as we discussed earlier is defined to be the dot-product between the corresponding feature vectors. For more details see [8, 9]

## 6.2 The Count-Min sketch

The Count-Min sketch is the simplest sketch and has found many applications. It was introduced by G. Cormode, S. Muthukrishnan[7]. The algorithm takes an error bound  $\epsilon$  and an error probability bound  $\delta$  and provides an additive approximation guarantee. We begin with a simpler version.

Let  $h : [n] \rightarrow [1, k]$  be a function chosen uniformly from a 2-universal family of hash functions. A 2-universal family  $\mathcal{H}$  of hash functions have that property,

$$\Pr_{h \in R\mathcal{H}}[h(x) = h(y)] = \frac{1}{k}$$

We keep an array  $\text{Count}[1 \dots k]$  of  $k$  integers. Consider the following simple algorithm.

---

**Algorithm** : Count-Min Sketch ( $k, \epsilon, \delta$ )

---

COUNT  $\leftarrow \text{ZEROS}(k)$  ▷ 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  $h(a_i)$

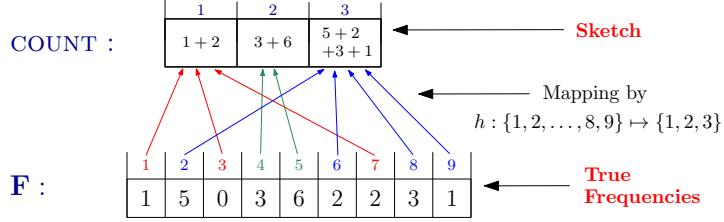
On query  $j$

▷ query:  $\mathbf{F}[j] = ?$

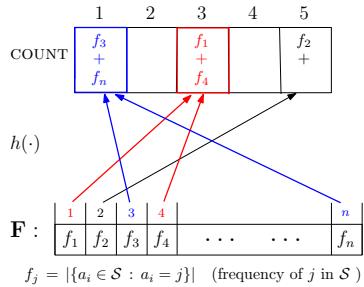
**return** COUNT[ $h(j)$ ]

---

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



Note that when  $k < n$  (which is typically the case, actually we require that  $k \in o(n)$  i.e.  $k \ll n$ ), the algorithm provides an upper bound on actual frequency (since the algorithm returns  $Count[h(j)]$ , other elements that hash to the same value, i.e. elements  $i$ , such that  $h(i) = h(j)$  also contribute to the returned value  $Count[h(j)]$ ).



Let  $\tilde{f}_j = Count[h(j)]$  be the estimate provided by the algorithm for query  $j$  then from the above reasoning we get that

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

For  $j \in [n]$ , we estimate the excess (error),  $X_j$  in  $\tilde{f}_j$ . Clearly,  $X_j = \tilde{f}_j - f_j$ . Let  $1_{h(i)=h(j)}$  be the indicator random variable for the event  $h(i) = h(j)$ .

$$1_{h(i)=h(j)} = \begin{cases} 1 & \text{if } h(i) = h(j) \\ 0 & \text{otherwise} \end{cases}$$

Note that  $i$  makes contribution to  $\tilde{f}_j$  iff  $h(i) = h(j)$  ( $1_{h(i)=h(j)} = 1$ ) and when it does contribute, its contribution is exactly  $f_i$ . We therefore get that

$$X_j = \sum_{i \in [n] \setminus j} f_i \cdot 1_{h(i)=h(j)}.$$

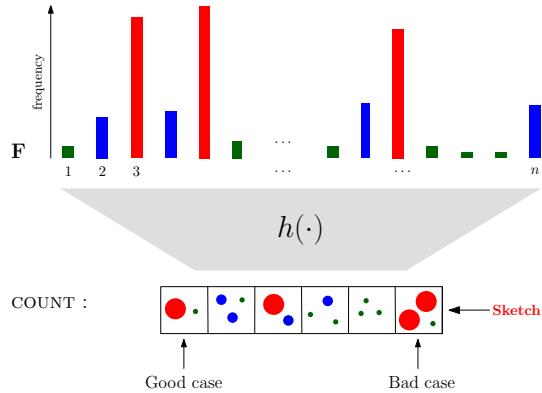
By the *goodness* of  $h$ , ( $h$  is a 2-universal hash function), we have that  $\forall i \neq j \quad \mathbb{P}[h(i) = h(j)] = \frac{1}{k}$ . This gives us

$$\mathbb{E}[1_{h(i)=h(j)}] = \mathbb{P}[h(i) = h(j)] = \frac{1}{k}.$$

We find the expectation of  $X_j$ .

$$\begin{aligned}
\mathbb{E}(X_j) &= \mathbb{E}\left(\sum_{i \in [n] \setminus j} f_i \cdot 1_{h(i)=h(j)}\right) \\
&= \sum_{i \in [n] \setminus j} \mathbb{E}(f_i \cdot 1_{h(i)=h(j)}) && \text{Linearity of expectation} \\
&= \sum_{i \in [n] \setminus j} f_i \cdot \mathbb{E}(1_{h(i)=h(j)}) && f_i \text{ is a constant} \\
&= \sum_{i \in [n] \setminus j} f_i \cdot \frac{1}{k} \\
&\leq \sum_{i \in [n] \setminus j} \|F\|_1 \cdot \frac{1}{k}
\end{aligned}$$

Where  $\|F\|_1$  is the sum of frequencies of all elements. Since all frequencies are non-negative it is actually the  $L_1$  norm of the frequency vector, that is why we denoted it by the  $L_1$  norm of  $F$ .



**Theorem 4** (Markov's Inequality). *If  $Z$  be a non-negative random variable, then*

$$\mathbb{P}[Z \geq t \cdot \mathbb{E}(Z)] \leq \frac{1}{t}$$

Substitute  $k = \frac{2}{\epsilon}$  (since  $k$ , the length of hash table is in our control) and using Markov's inequality we get that

$$\mathbb{P}[X_j \geq \epsilon \|F\|_1] = \mathbb{P}[X_j \geq 2 \mathbb{E}[X_j]] \leq \frac{1}{2}$$

Summarizing the analysis of this algorithm we get

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

Hence Algorithm 1 is an  $(\epsilon \|F\|_1, \frac{1}{2})$ -additive approximation algorithm. Space required by the algorithm is  $k$  integers (plus some more for processing etc.) and  $k = \frac{2}{\epsilon}$  is a constant.

### 6.2.1 Amplifying the probability

We can amplify the probability of success by selecting  $t$  independent hash functions  $h_1, h_2, \dots, h_t$  each from a 2-universal family of hash functions and proceed as follows. Each  $h_i : [n] \rightarrow [1 \dots k]$

---

**Algorithm** : Count-Min Sketch ( $k, \epsilon, \delta$ )

---

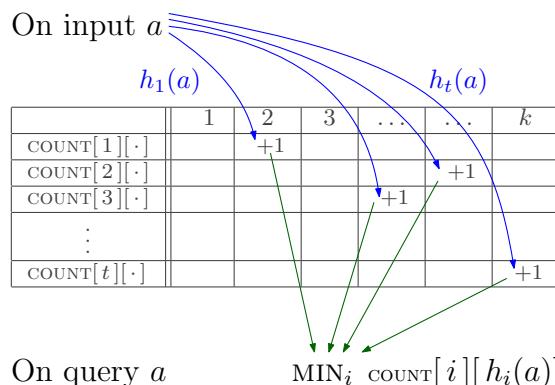
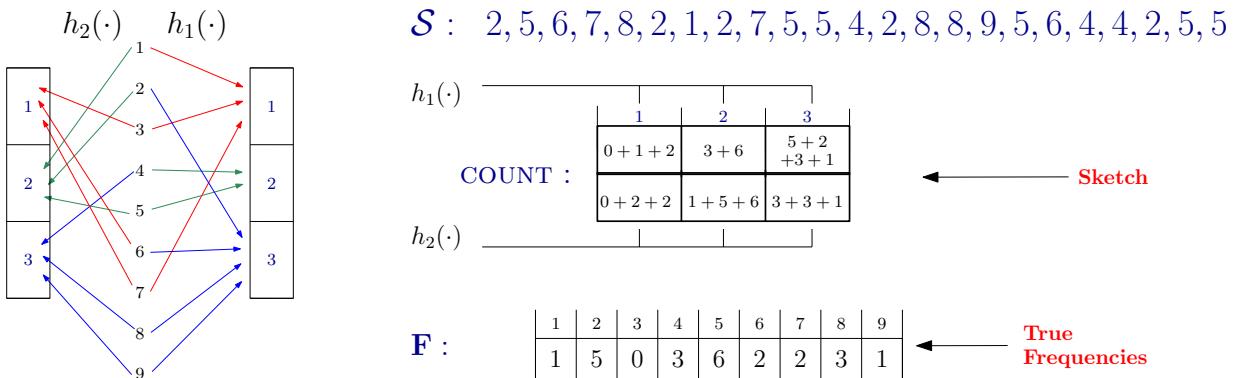
```

COUNT  $\leftarrow$  ZEROS( $t \times k$ )                                 $\triangleright$  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
                                             $\triangleright$  increment COUNT[ $r$ ] at index  $h_r(a_i)$ 
On query  $j$                                  $\triangleright$  query:  $\mathbf{F}[j] = ?$ 
return  $\min_{1 \leq r \leq t} \text{COUNT}[r][h_r(j)]$ 

```

---

So we keep  $t$  estimates instead of 1, and since every estimate is an upper bound, it is clear that we should return the minimum of all estimates (the one which has minimum contribution from other elements).



Define  $X_{jr}$  to be the contribution of other elements to  $\text{Count}[r][h(j)]$ . We know that

$$\mathbb{P}[X_{jr} \geq \epsilon \|F\|_1] \leq \frac{1}{2}$$

if the length of hash table is  $k = \frac{2}{\epsilon}$ .

Now if  $\tilde{f}_j \geq f_j + \epsilon \|F\|_1$ , then for all  $1 \leq r \leq t$ , we must have that  $X_{jr} \geq \epsilon \|F\|_1$ , and the probability of this event is (since  $h_r$ 's are independent) is bounded as

$$\mathbb{P}[\forall r X_{jr} \geq \epsilon \|F\|_1] \leq \left(\frac{1}{2}\right)^t$$

Substitute  $t = \log(\frac{1}{\delta})$  (the number of hash functions is in our control) and we get

$$\mathbb{P}[\forall r X_{jr} \geq \epsilon \|F\|_1] \leq \left(\frac{1}{2}\right)^{\log 1/\delta} = \delta$$

Summarizing the analysis of this algorithm 2 we get

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

Hence Algorithm 2 is an  $(\epsilon \|F\|_1, \delta)$  additive approximation algorithm. Space required by the algorithm is  $k \cdot t$  integers (plus some more for processing etc), and  $k = \frac{2}{\epsilon}$  and  $t = \log(1/\delta)$ .

## 7 Finding frequency of an element: The Count sketch

Since the Count-Min sketch always provides an upper bound, errors always accumulate, the Count sketch attempts to have the errors cancel each other. It was introduced [6]. Again we begin with a simpler version and later extend it.

### 7.1 The Count sketch

Let  $h : [n] \rightarrow [1, k]$  and  $g : [n] \rightarrow \{+1, -1\}$  be two hash functions independent of each other. We keep an array  $\text{Count}[1 \dots k]$  of  $k$  integers. Consider the following simple algorithm.

---

**Algorithm** : Count Sketch  $(k, \epsilon, \delta)$

---

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  $a_i$

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

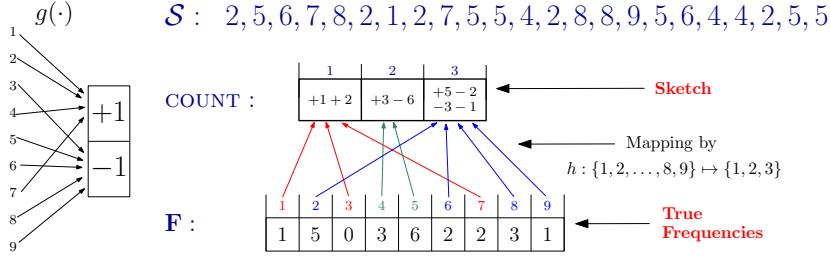
▷ increment or decrement, depending on value of  $g(a_i)$  COUNT at index  $h(a_i)$

On query  $j$

▷ query:  $\mathbf{F}[j] = ?$

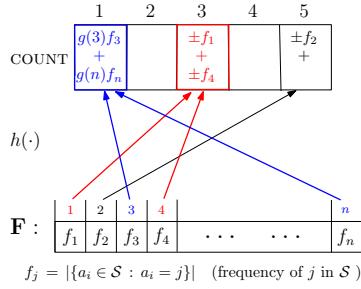
**return**  $g(j) \times \text{COUNT}[h(j)]$

---



Let  $\tilde{f}_j$  be the estimate provided by the algorithm for query  $j$ , (i.e.  $\tilde{f}_j = g(j) \cdot \text{Count}[h(j)]$ ), we will first show that in expectation the algorithm provides the correct answer

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



For  $j \in [n]$ , we will estimate the error in  $\tilde{f}_j$ . Examining the algorithm it is easy to see that the  $i \in [n]$  contributes to  $\text{Count}[h(j)]$  iff  $h(i) = h(j)$  and when it does contribute, its contribution is either  $f_i$  or  $-f_i$ , depending on  $g(i)$ .

From the above reasoning we get  $\text{Count}[h(j)] = \sum_{i \in [n]} f_i \cdot g(i) \cdot 1_{h(i)=h(j)}$ . On query  $j$  the output of the algorithm,  $\tilde{f}_j$

$$\begin{aligned} \tilde{f}_j &= g(j) \cdot \text{Count}[h(j)] = g(j) \cdot \sum_{i \in [n]} f_i \cdot g(i) \cdot 1_{h(i)=h(j)} = g(j) \cdot \left( f(j) \cdot g(j) + \sum_{i \in [n] \setminus j} f_i \cdot g(i) \cdot 1_{h(i)=h(j)} \right) \\ &= f(j) \cdot (g(j))^2 + \sum_{i \in [n] \setminus j} f_i \cdot g(i)g(j) \cdot 1_{h(i)=h(j)} = f(j) + \sum_{i \in [n] \setminus j} f_i \cdot g(i)g(j) \cdot 1_{h(i)=h(j)} \end{aligned}$$

The expected value of  $\tilde{f}_j$  is

$$\begin{aligned}
\mathbb{E}(\tilde{f}_j) &= \mathbb{E} \left( f(j) + \sum_{i \in [n] \setminus j} f_i \cdot g(i)g(j) \cdot 1_{h(i)=h(j)} \right) \\
&= f(j) + \mathbb{E} \left( \sum_{i \in [n] \setminus j} f_i \cdot g(i)g(j) \cdot 1_{h(i)=h(j)} \right) && f_j \text{ is a constant} \\
&= f(j) + \sum_{i \in [n] \setminus j} \mathbb{E} (f_i \cdot g(i)g(j) \cdot 1_{h(i)=h(j)}) && \text{linearity of expectation} \\
&= f(j) + \sum_{i \in [n] \setminus j} f_i \cdot \mathbb{E} (g(i)g(j)) \cdot \mathbb{E} (1_{h(i)=h(j)}) && h \text{ and } g \text{ are independent}
\end{aligned}$$

As above  $\mathbb{E}(1_{h(i)=h(j)}) = \mathbb{P}[h(i) = h(j)] = \frac{1}{k}$  and  $\mathbb{E}(g(i)g(j)) = -1(\frac{1}{2}) + 1(\frac{1}{2}) = 0$ . Hence we get that

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

Next we calculate the variance of  $\tilde{f}_j$

$$\begin{aligned}
Var(\tilde{f}_j) &= Var \left( f(j) + \sum_{i \in [n] \setminus j} f_i \cdot g(i)g(j) \cdot 1_{h(i)=h(j)} \right) \\
&= Var \left( \sum_{i \in [n] \setminus j} f_i \cdot g(j)g(i) \cdot 1_{h(i)=h(j)} \right) && Var(a+X) = Var(X) \\
&= (g(j))^2 \cdot Var \left( \sum_{i \in [n] \setminus j} f_i \cdot g(i) \cdot 1_{h(i)=h(j)} \right) && Var(aX) = a^2 Var(X) \\
&= Var \left( \sum_{i \in [n] \setminus j} f_i \cdot g(i) \cdot 1_{h(i)=h(j)} \right) && (g(j))^2 = 1 \\
&= \sum_{i \in [n] \setminus j} Var(f_i \cdot g(i) \cdot 1_{h(i)=h(j)}) && g \text{ and } h \text{ are independent} \\
&= \sum_{i \in [n] \setminus j} \mathbb{E}(f_i \cdot g(i) \cdot 1_{h(i)=h(j)})^2 - (\mathbb{E}(f_i \cdot g(i) \cdot 1_{h(i)=h(j)}))^2 && Var(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 \\
&= \sum_{i \in [n] \setminus j} \mathbb{E}(f_i \cdot g(i) \cdot 1_{h(i)=h(j)})^2 - \sum_{i \in [n] \setminus j} (\mathbb{E}(f_i \cdot g(i) \cdot 1_{h(i)=h(j)}))^2 && \text{Linearity of Exp.} \\
&= \sum_{i \in [n] \setminus j} \mathbb{E}(f_i \cdot g(i) \cdot 1_{h(i)=h(j)})^2 - 0 && \text{from above} \\
&= \sum_{i \in [n] \setminus j} f_i^2 \mathbb{E}(g(i) \cdot 1_{h(i)=h(j)})^2 - 0 && \mathbb{E}(aX) = a \mathbb{E}(X) \\
&= \sum_{i \in [n] \setminus j} f_i^2 \mathbb{E}(g(i)^2 \cdot 1_{h(i)=h(j)}^2) = \sum_{i \in [n] \setminus j} f_i^2 \cdot \frac{1}{k} \leq \sum_{i \in [n]} f_i^2 \cdot \frac{1}{k} \\
&\leq \sum_{i \in [n]} (\|F\|_2)^2 \cdot \frac{1}{k}
\end{aligned}$$

where  $\|F\|_2$  is the  $L_2$  norm of the frequency vector.

**Theorem 5** (Chebyshev's Inequality). *If  $Z$  be a random variable, then*

$$\mathbb{P}[|Z - \mathbb{E}(Z)| \geq t\sqrt{Var(X)}] \leq \frac{1}{t^2}$$

Apply the Chebyshev's inequality with  $Var(\tilde{f}_j) \leq \frac{(\|F\|_2)^2}{k}$ , we get that

$$\mathbb{P}\left[|\tilde{f}_j - f_j| \geq \epsilon \sqrt{k} \sqrt{\frac{(\|F\|_2)^2}{k}}\right] \leq \frac{1}{k\epsilon^2} \implies \mathbb{P}\left[|\tilde{f}_j - f_j| \geq \epsilon \|F\|_2\right] \leq \frac{1}{k\epsilon^2}$$

Substitute  $k = \frac{3}{\epsilon^2}$  (since  $k$  is length of the hash table which in our control) and using Chebychev's inequality we get that

$$\mathbb{P} \left[ |\tilde{f}_j - f_j| \geq \epsilon \|F\|_2 \right] \leq \frac{1}{3}$$

Summarizing the analysis of this algorithm we get

- $\mathbb{E}(\tilde{f}_j) = f_j$
- $\tilde{f}_j \leq f_j + \epsilon \|F\|_2$  with probability at least  $\frac{1}{3}$
- $\tilde{f}_j \geq f_j - \epsilon \|F\|_2$  with probability at least  $\frac{1}{3}$  OR
- $f_j - \epsilon \|F\|_2 \leq \tilde{f}_j \leq f_j + \epsilon \|F\|_2$  with probability at least  $\frac{1}{3}$

Hence Algorithm 3 is an  $(\epsilon \|F\|_2, \frac{1}{3})$ -additive approximation algorithm. Space required by the algorithm is  $k$  integers (plus some more for processing etc), and  $k = \frac{3}{\epsilon^2}$  is a constant.

### 7.1.1 Amplifying the probability

We can amplify the probability of success by selecting  $t$  independent hash functions  $h_1, h_2, \dots, h_t$  and  $t$  independent hash functions (signs)  $g_1, g_2, \dots, g_t$  each from a 2-universal family of hash functions and proceed as follows. Each  $h_i : [n] \rightarrow [1 \dots k]$  and  $g_i : [n] \rightarrow \{+1, -1\}$

---

**Algorithm** : Count Sketch  $(k, \epsilon, \delta)$

---

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:  $\mathbf{F}[j] = ?$

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

---

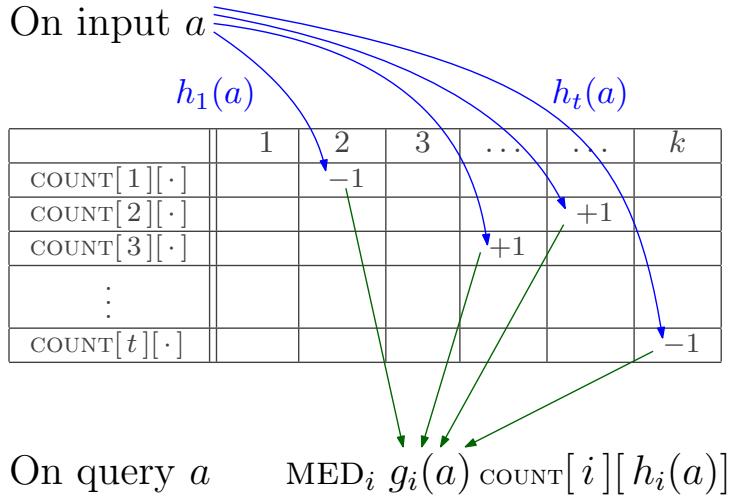
It is clear that why should we return median of the  $t$  estimates, as each one of them can be an over or under estimate.

For the analysis, let  $\tilde{f}_{jr}$ , ( $1 \leq r \leq t$ ), be the  $r$ -th estimate, i.e.

$$\tilde{f}_{jr} = g_r(j) \cdot \text{Count}[r][h_r(j)]$$

From above analysis we know that when  $k = \frac{3}{\epsilon^2}$ ,

$$\mathbb{P} \left[ |\tilde{f}_{jr} - f_j| \geq \epsilon \|F\|_2 \right] \leq \frac{1}{3}$$



Note that since we returned the median of  $t$  estimates, if the median is more than  $x$  units away from the actual value, then at least half of all the estimates are more than  $x$  units away from the actual value. The probability of many independent estimates being far away is much smaller. More precisely, Let

$$X_r = \begin{cases} 1 & \text{if } |\tilde{f}_{jr} - f_j| \geq \epsilon \|F\|_2 \\ 0 & \text{otherwise} \end{cases}$$

Basically  $X_r$  is the indicator random variable, indicating whether  $\tilde{f}_{jr}$  has too much error. We count how many of the estimates have too much error, i.e. let

$$X = \sum_{r=1}^t X_r.$$

Now since we said that if our answer is wrong (has error beyond the  $\epsilon \|F\|_2$  margin), then at least half of the estimates would be wrong. In other words  $X$  would be more than  $t/2$ . But we will show that only with very small probability  $X$  could be so large. By the way we have

$$\mathbb{E}(X) = \sum_{r=1}^t \mathbb{E}(X_r) = \sum_{r=1}^t \frac{1}{3} = \frac{t}{3}$$

**Theorem 6** (Chernoff Bound). *Let  $Z_1, Z_2, \dots, Z_n$  be independent random variables,  $0 \leq Z_i \leq 1$ . Let  $Z = Z_1 + Z_2 + \dots + Z_n$ . Suppose  $\mu = \mathbb{E}(Z) = \mathbb{E}(Z_1) + \mathbb{E}(Z_2) + \dots + \mathbb{E}(Z_n)$ . Then for any  $\gamma \geq 0$*

- $\mathbb{P}[Z \geq (1 + \gamma)\mu] \leq \exp\left(-\frac{\gamma^2}{2 + \gamma}\mu\right)$
- $\mathbb{P}[Z \leq (1 - \gamma)\mu] \leq \exp\left(-\frac{\gamma^2}{2}\mu\right)$

The computer science version of this says that  $\mathbb{P}[|Z - \mu| \geq \gamma\mu] \leq 2^{-\Omega(\mu\gamma^2)}$

Since the expected number of erroneous estimates is  $t/3$ , if the median of estimates is wrong, then the number of erroneous estimates must be at least  $t/2 = t/3 + t/6$ , i.e.  $X$  must deviate from its expected value by at least  $t/3$  and by the Chernoff bound the probability of that happening is given as follows:

$\mathbb{P}[X \geq (1 + 1/2)\frac{t}{3}] \leq \exp\left(-\frac{(1/2)^2}{2 + (1/2)} \cdot \frac{t}{3}\right) = \exp\left(-\frac{t}{30}\right)$ . We want this probability to be at most  $\delta$  (the given bound on error probability), substitute  $t = 30 \ln(\frac{1}{\delta})$ . Note that  $t$  is in our control, the number of times we want to run Algorithm 3 in parallel. Overall we get,

- $\mathbb{E}(\tilde{f}_j) = f_j$
- $|\tilde{f}_j - f_j| \leq \epsilon \|F\|_2$  with probability at least  $1 - \delta$

Hence Algorithm 4 is an  $(\epsilon \|F\|_2, \delta)$  additive approximation algorithm. Space required by the algorithm is  $k \cdot t$  integers (plus some more for processing etc), and  $k = \frac{3}{\epsilon^2}$  and  $t = 30 \ln(1/\delta)$  (assuming my calculations are right).

## 8 Comparison of Count and Count-Min Sketch

Note that the Count sketch uses more space than the Count-Min sketch, but does it give better estimates. Recall our discussion on the  $L_k$  norms.

- The deviation for the Count-Min sketch is bounded by  $\epsilon \|F\|_1$  and that of the Count sketch is bounded by  $\|F\|_2$ .
- For all vectors  $z \in \mathbb{R}^n$ ,  $\|z\|_1 \geq \|z\|_2$ , hence Count-Min sketch will always give looser bound.
- The two give exactly the same answer when only one element has non-zero frequency, since for  $F = (f_1, 0, 0, 0, 0, 0)$ ,  $\|F\|_1 = \|F\|_2 = 1$ .
- Count outperforms the Count-Min sketch by the most when all frequencies are equal. Say  $F = (1, 1, 1, \dots, 1)$ , in this case  $\|F\|_1 = n$  and  $\|F\|_2 = \sqrt{n}$ . So for say  $\epsilon = 1/10$ , the Count-Min estimate is within  $n/10$  of the actual answer (with probability  $1 - \delta$ ), which the Count estimate is within  $\sqrt{n}/10$  of the actual answer.
- As the distribution of frequencies spreads out, Count outperforms the Count-Min.

## 9 Lower bound

In this section, I want to show that the Count-Min sketch is asymptotically tight with respect to space requirements. This is to give you a general idea of how lower bounds work and what do they mean. Please read it on your own and if you encounter something like this in your research, you would be familiar with the typical structure and terminology. We use a reduction from the well known Index problem from communication complexity.

## 9.1 Index problem

In this problem there are two players, Alice and Bob. Alice holds a binary array  $A$  of length  $n$  (bit string). Bob has an index  $j \in [n]$ . Bob wants to find out the value of  $A[j]$ . All communication has to be one-way, i.e. only Alice can send information to Bob. Bob knows the value of  $n$  but doesn't know the contents of  $A$ . We denote an instance of index problem as the pair  $(A, j)$

The question is how many bits Alice needs to send to Bob so as he determines the value of  $A[j]$ . We want to reduce the communication complexity of the protocol (algorithm). We only require that Bob determines  $A[j]$  with high probability (one-way probabilistic communication).

We discussed that

- It is very easy if Alice is allowed to send all  $n$  bits, as then Bob will check  $A[j]$ . This requires  $n$  bits communication.
- With some XoR tricks, maybe we can reduce to  $n - 1$  or so bits communication.
- Bob can randomly guess the value of  $A[j]$ , and his guess will be correct with probability  $1/2$ . But we want it to be correct with probability  $1 - \delta$  for a small input parameter  $\delta$ .
- We discussed why would one be interested in Index problem, Private Information Retrieval etc.

We stated the following theorem from information theory.

**Theorem 7.** (*Lower bound on index problem*) *For Bob to determine the  $A[j]$  correctly with probability at least  $1/2 - \delta$ , Alice needs to send  $\Omega(\delta^2 n)$  bits.*

So basically it says that to improve upon the random guess the communication needed is  $\Omega(n)$  bits.

## 9.2 Reduction from the Frequency estimation to Index problem

**Theorem 8.** *To estimate  $f_j$  within error  $\epsilon \|F\|_1$  with constant probability, one needs at least  $\Omega(1/\epsilon)$  space.*

*Proof.* Assume that there is a solution that estimates  $f_j$  within error  $\epsilon \|F\|_1$  with probability at least  $1 - \delta$  and it uses  $O(1/\epsilon)$  space. This means there is an algorithm  $G$  that outputs estimate  $\tilde{f}_j$  such that  $\mathbb{P}[\lvert \tilde{f}_j - f_j \rvert \leq \epsilon \|F\|_1] \leq 1 - \delta$  and  $G$  uses space at most  $1/\epsilon$ .

We will show that if there is such an algorithm  $G$ , then using  $G$  we can efficiently solve the Index problem. To see that suppose we have an instance of the index problem  $(A, j)$ . We make a stream  $\mathcal{S}$  with each element from the set  $\{0, 1, \dots, n\}$ . The stream will realize the following frequency vector.

$f_0 = 2 \cdot |\{i : A[i] = 0\}|$ , i.e. the frequency of 0 is twice the number of 0's in  $A$ .  
 $f_i = 2$ , if  $A[i] = 1$ .

In other words, for each 0 in  $A$  we put two 0's in  $\mathcal{S}$  and for each 1 at location  $i$ , we put two 1's in  $\mathcal{S}$ . As an example, for the following  $A$

1	2	3	4	5	6
1	0	1	0	1	1
0	1	2	3	4	5
4	2	0	2	0	2

The stream  $\mathcal{S}$  will be 1, 1, 0, 0, 3, 3, 0, 0, 5, 5, 6, 6

It is easy to see that this  $\mathcal{S}$  defines the following frequency vector  $F$

Let  $\epsilon = 1/2n$ , clearly  $\|F\|_1 = 2n = 1/\epsilon$  and  $\epsilon\|F\|_1 = 1$ . Lets apply the algorithm  $G$  on  $\mathcal{S}$  and estimate the frequency of  $j$  (within error  $\epsilon\|F\|_1$ ).

If for  $\tilde{f}_j$ ,  $G$  outputs anything above 1, then  $A[j] = 1$  and otherwise  $A[j] = 0$ . Because if  $A[j] = 0$ , then  $f_j = 0$ , and the algorithm can overestimate it by at most  $0+1 = \epsilon\|F\|_1$  with high probability  $G$ , and if  $A[j] = 1$ , then  $f_j = 1$  and the algorithm will return something between 1 and 3 with high probability.

The contradiction comes from the fact that now Alice can make a stream  $\mathcal{S}$  given the array  $A$  with her. And run algorithm  $G$  on  $\mathcal{S}$  and send to Bob whatever data structure  $G$  is using (the sketch  $G$  computes). Since it was claimed that  $G$  uses space  $O(1/\epsilon)$ , hence the communication from Alice to Bob is at most  $O(1/\epsilon) = O(n)$ , which contradicts the theorem above.

□

## References

- [1] S. Ali, H. Mansoor, N. Arshad, and I. Khan. Short term load forecasting using smart meter data. In *International Conference on Future Energy Systems*, pages 419–421, 2019.
- [2] S. Ali, H. Mansoor, I. Khan, N. Arshad, S. Faizullah, and M.A. Khan. Fair allocation based soft load shedding. In *Proceedings of SAI Intelligent Systems Conference*, pages 407–424, 2020.
- [3] S. Ali, H. Mansoor, I. Khan, N. Arshad, Muhammad A. Khan, and S. Faizullah. Short-term load forecasting using ami data. *arXiv preprint arXiv:1912.12479*, 2020.
- [4] Sarwan Ali, Maria Khalid Alvi, Safi Faizullah, Muhammad Asad Khan, Abdullah Alshanti, and Imdadullah Khan. Detecting ddos attack on sdn due to vulnerabilities in openflow. In *2019 International Conference on Advances in the Emerging Computing Technologies (AECT)*, pages 1–6, 2020.
- [5] M. Beg, M. Ahmad, A. Zaman, and I. Khan. Scalable approximation algorithm for graph summarization. In *Advances in Knowledge Discovery and Data Mining (PAKDD)*, pages 502–514, 2018.
- [6] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *Proceedings of the 29th International Colloquium on Automata, Languages and Programming, ICALP '02*, page 693703. Springer-Verlag, 2002.

- [7] Graham Cormode and S. Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *Journal of Algorithms*, 55(1):58–75, 2005.
- [8] M. Farhan, J. Tariq, A. Zaman, M. Shabbir, and I. Khan. Efficient approximation algorithms for strings kernel based sequence classification. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 6935–6945, 2017.
- [9] P. Kuksa, I. Khan, and V. Pavlovic. Generalized similarity kernels for efficient sequence classification. In *SIAM Intern. Conf. on Data Mining (SDM)*, pages 873–882, 2012.