

# Algorithms for streaming data

Feb 8 2022

- Limited storage space
- data streaming past

Notation:

Denote elements of the stream

$a_1, a_2, a_3, \dots$

$a_i \in U$

Functions:

1. sum of all elements
2. max
3. heavy-hitters
4. median
5. # unique elements

Sampling?

Can lead to incorrect answers (if not done correctly)

Example: Find "# uniques" = # elements that occur exactly once.

Sampling approach (can lead to incorrect answers)

- Sample 10% of the stream picking each element of the stream w.p. 0.1

- Count # uniques among sampled elements

- Multiply by 10

Can lead to incorrect answers. Example ...

Suppose stream length =  $n$   
 $n/2$  uniques  
 $n/4$  appear twice

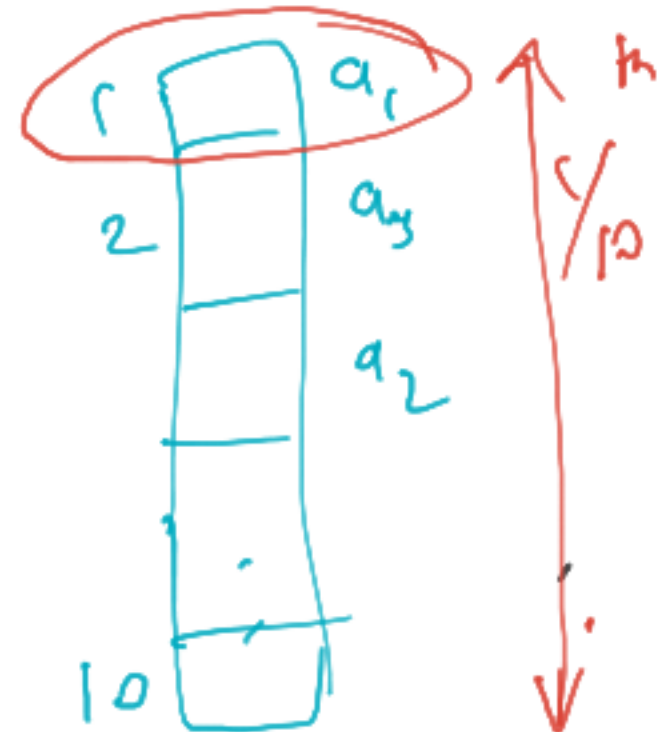
In the sampled stream,  
 Expected length =  $\frac{n}{10}$

$$\# \text{ uniques} = \frac{n}{2} \times 0.1 + \frac{n}{4} (2 \times 0.1 - (0.1)^2)$$

$$\approx \frac{n}{10}$$

$$\text{Estimate} \approx \frac{n}{10} \times 10 = \underline{\underline{n}}$$

$a_1, a_2, a_1$   
 $\uparrow \quad \uparrow$   
 $\times$



Abstraction: view stream as a vector (in high dim. space)

Stream at time  $t$   $x^t \in \mathbb{Z}^{|U|}$

$$x^t = (x_1^t, x_2^t, \dots, x_{|U|}^t)$$

vector of  
length  
 $|U|$

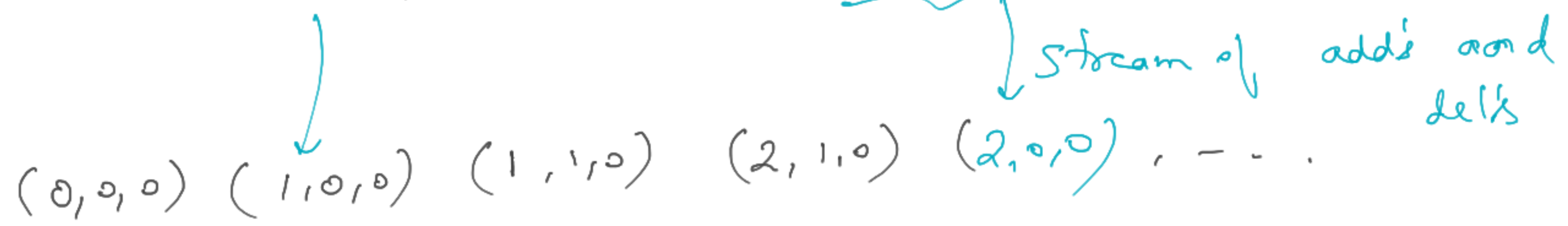
$x_i^t =$  num. of times the  $i^{\text{th}}$  element of  $U$   
has been seen until time  $t$

E.g. if  $a_j \in U$  appears,  $x_j + 1$

- Allow both additions and deletions

E.g.  $U = \{A, B, c\}$

add(A), add(B), add(A), del(B), del(A), add(c)



Heavy hitters: estimate "large" entries

Total num. of element seen:

# distinct elements = # non-zero entries

# Heavy hitters:

$\epsilon$ -heavy hitters: Indices  $i$  s.t.

$$x_i > \epsilon \|x\|_1$$

$\|\cdot\|_1 = \text{sum of the elements}$

Simpler problem first: "Count-query"

At any time  $t$ , any given index  $i$ ,  
find an estimate of  $x_i^t$  with  
error at most  $\epsilon \|x\|_1$

i.e.  $\overset{t}{x_i} \in x_i \pm \epsilon \|x\|_1$

estimate  $\rightarrow$

To the first order, we can look at  $i$ 's s.t.  $\overset{t}{x_i} > 0$

Hashing-based solution : Count-min Sketch

2-Step soln.

Step 1: Let  $h: U \rightarrow [M]$  hash fn.

Let  $A[1 \dots M]$  array for storing non-negative integers

When update  $a_t$  arrives

If  $(a_t == \text{add}, i)$

then  $A[h(a_t)]++$

else (i.e.  $\text{del}, i$ )

$A[h(a_t)]--$

Estimate for  $x_i^t$  :  $\hat{x}_i^t = A[h(i)]$

$$\tilde{x}_i^t = A[h(i)]$$

$$= x_i^t +$$

$$\sum_{\substack{j \in U \\ j \neq i}} x_j^t \underbrace{\mathbb{1}[h(j) = h(i)]}_{\text{Collisions at index } h(i)}$$

added later: (over estimate)

$$\text{Error} = \tilde{x}_i^t - x_i^t =$$

$$\mathbb{E}[\text{error}] = \sum_{\substack{j \in U \\ j \neq i}} x_j^t \mathbb{E}[\mathbb{1}[h(j) = h(i)]]$$

$$= \sum_{\substack{j \in U \\ j \neq i}} x_j^t P[h(j) = h(i)]$$

Assume "universal hashing"  $\Rightarrow \leq \sum_{j \neq i} x_j^t \cdot \frac{1}{M}$



$E[\text{error}]$

$$\leq \frac{\sum_{j \neq i} x_j^t}{M}$$
$$\leq \frac{\|x^t\|_1}{M}$$

$$\varepsilon \|x^t\|_1$$
$$\varepsilon = \frac{1}{M}$$

In expectation.

Can we give probabilistic guarantee?

$$P(|\bar{x}_i^t - x_i^t| < \varepsilon \|x^t\|_1) \geq 1 - \delta$$

Step 2:

$l$  hash functions  $h_1, h_2, \dots, h_l : U \rightarrow [M]$

$l$  arrays  $A_1, \dots, A_l$

E.g.:

$i$ th element arrive:  $A_k[h_k(i)]++$  for  $k=1, \dots, l$

overall estimate  $x_i = \min_{k=1}^l A_k[h_k(i)]$

For any single estimator,

$$P \left[ \text{error} > \underbrace{2 \cdot \frac{\|x^t\|}{M}}_{\alpha} \right] \leq \frac{1}{2}$$

Markov Inequality

Non. neg. R.v.  $X$   
with mean  $\mu$

$$P(X \geq \underline{\alpha}) \leq \frac{\mu}{\alpha}$$

Now,

$$P \left[ \text{error of } x_i^t > 2 \cdot \frac{\|x^t\|_1}{M} \right]$$

$$= \left( \frac{1}{2} \right)^\lambda$$

$$P \left[ \|x_i^t - x_i\| \leq \epsilon \|x^t\|_1 \right] \geq 1 - \delta$$

where  $M = \frac{2}{\epsilon}$ ,  $\delta = \left( \frac{1}{2} \right)^\lambda$

$$\Rightarrow \lambda = \log_2 \left( \frac{1}{\delta} \right)$$

Tells us how to choose the parameters.