

# CS7800: Advanced Algorithms

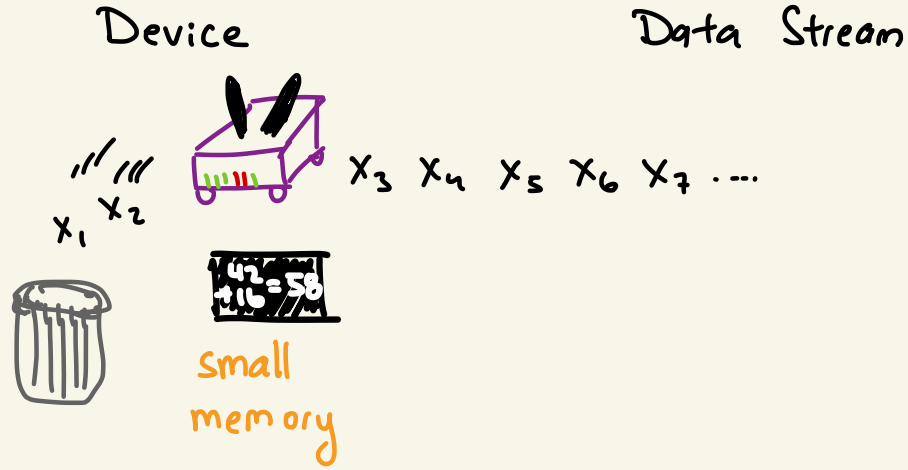
## Lecture 20: Streaming I

- Streaming Algorithms
- Distinct Elements

Jonathan Ullman

11-29-2022

# Streaming Algorithms



- Stream of elements  $x_1, x_2, x_3, \dots$  from a universe  $U$
- Want to only store a small number of bits  $S$  at any time
- Want to approximate some function  $f(x_1, x_2, x_3, \dots)$

# Warmup: Uniform Sampling in a Stream

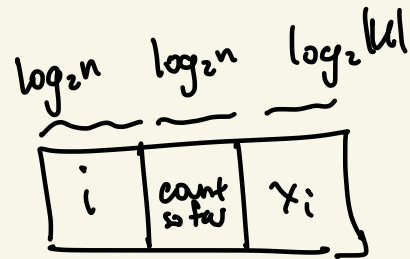
Inputs: a stream of elements  $x_1, x_2, x_3, \dots$  from  $U$

Goal: output a uniformly random element  $s$  from the stream

↑  
If stream has length  $n$   
then  $P(s=x_i) = \frac{1}{n}$

If you know the stream length  $n$ :

- Pick  $i \in \{1, 2, \dots, n\}$  at random
- Wait until you see  $x_i$
- Write down  $x_i$



$$O(\log n + \log(|U|))$$

# Reservoir Sampling

$x_1$   $x_2$   $x_3$   $x_4$   $x_5 \dots x_i$

$$s = x_1$$

$$s = \frac{1}{2}x_1 + \frac{1}{2}x_2$$

w.p.  $\frac{1}{2}$ , let  $s = x_2$   
w.p.  $\frac{1}{2}$ , do nothing

w.p.  $\frac{1}{i}$ , let  $s = x_i$   
w.p.  $1 - \frac{1}{i}$ , do nothing

$$s = \frac{2}{3}\left(\frac{1}{2}x_1 + \frac{1}{2}x_2\right) + \frac{1}{3}x_3$$
$$= \frac{1}{3}x_1 + \frac{1}{3}x_2 + \frac{1}{3}x_3$$

w.p.  $\frac{1}{3}$ , let  $s = x_3$   
w.p.  $\frac{2}{3}$ , do nothing

space is  $O(\log n + \log |U|)$

# Reservoir Sampling

for  $i=1, 2, 3, \dots$ :

get  $x_i$

w.p.  $\frac{1}{i}$  set  $s_i = x_i$

w.p.  $1 - \frac{1}{i}$  set  $s_i = s_{i-1}$

} only store  $s_i$

return  $s_n$

Thm: For every  $n$ ,  $s_n$  is a random sample from  $\{x_1, x_2, \dots, x_n\}$

Pf: Base case is easy

Assume that  $s_{n-1}$  is random from  $\{x_1, \dots, x_{n-1}\}$   $\left( \begin{array}{l} P(s_{n-1} = x_i) = \frac{1}{n-1} \\ \text{for } i=1, \dots, n-1 \end{array} \right)$

$$P(s_n = x_n) = \frac{1}{n} \text{ (by construction)}$$

$$P(s_n = x_i) = \left(1 - \frac{1}{n}\right) \cdot \left(\frac{1}{n-1}\right) = \frac{1}{n} \quad \square$$

for  $i \leq n-1$

# Reservoir Sampling

for 2 samples w/o replacement

$x_1 \quad x_2 \quad \left. \vphantom{x_1} \right\} x_3 \quad \left. \vphantom{x_1} \right\} x_4 \quad x_5 \quad \dots \quad x_i \quad \left. \vphantom{x_1} \right\}$

$S = \{x_1, x_2\}$

$\frac{1}{3}(x_1, x_2)$   
 $+\frac{1}{3}(x_2, x_3)$   
 $+\frac{1}{3}(x_1, x_3)$

w.p.  $1 - \frac{2}{i}$ , hold  
w.p.  $\frac{2}{i}$ , keep  $x_i$   
and one of your  
samples at random

$$\frac{2(i-1)}{i(i-1)} = \frac{2}{i}$$

w.p.  $\frac{1}{3}$ , do nothing

w.p.  $\frac{2}{3}$ , pick one of  $S$  at random  
and replace the other with  $x_3$

# Reservoir Sampling

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	...
0.623	0.715	0.112	0.100	0.563	

$$h(x) = \text{md5}(s \parallel x)$$

# Counting Distinct Elements

Inputs: a stream of elements  $x_1, x_2, x_3, \dots$  from  $U$

Goal: the (approximate) number of distinct elements in the stream

$c$ -approximate means

$$\frac{1}{c} \cdot DE \leq \tilde{DE} \leq c \cdot DE$$

with high probability

$$DE_x = |\{u \in U : x_i = u \text{ for some } i\}|$$

stream: 1 1 3 4 8 3 1 2 8 3

#distinct: 5

Baseline: storing all the elements you've seen so far takes

$DE_x \cdot \log |U|$  bits of space

store a flag for each element takes  $|U|$  bits of space



# A Simplification: Threshold Testing

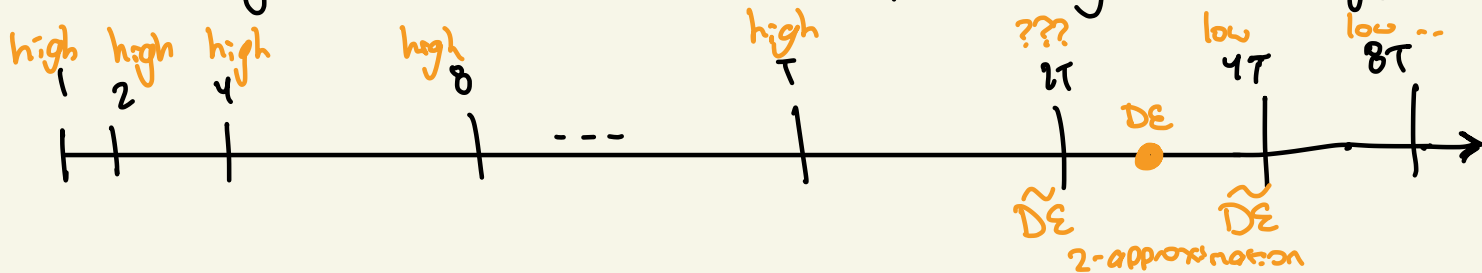
Goal is to design an algorithm  $A_T$  such that

① If  $DE \leq T$  then  $IP(A_T = \text{low}) \geq 1 - \delta$

② If  $DE \geq \underline{2T}$  then  $IP(A_T = \text{high}) \geq 1 - \delta$

2 is arbitrary

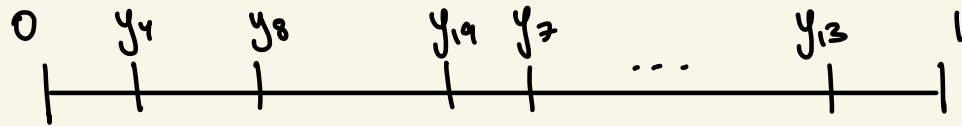
Can run  $\log_2 |I|$  copies of  $A_1, A_2, A_4, A_8, \dots, A_{|I|}$  in parallel to get a 2-approximation with probability  $\geq 1 - \delta \cdot \log_2 |I|$



# Threshold Testing Distinct Elements I

suppose the distinct elements are  $y_1, \dots, y_d$

suppose I assign a random number in  $[0, 1]$  to each



$$\approx \frac{1}{d+1}$$

use hashing to assign the same number to multiple copies of one elt

Information about the # of distinct elements is found in the widths of these gaps

# Threshold Testing Distinct Elements I

Choose a <sup>uniformly</sup> random hash function

$$h: U \rightarrow \{0, 1, 2, \dots, \underbrace{T-1}_{\text{threshold}}\}$$

For each  $x_i$  in the stream:

[ If  $h(x_i) = 0$  output high

Output low

Suppose  $D \leq T$

$$P(\text{low}) = \left(1 - \frac{1}{T}\right)^T \approx \frac{1}{e}$$

want  $P(\text{low}) \geq 1 - \delta$

Suppose  $D \geq 2T$

$$\begin{aligned} P(\text{low}) &= \left(1 - \frac{1}{T}\right)^{2T} \\ &= \left(\left(1 - \frac{1}{T}\right)^T\right)^2 \approx \frac{1}{e^2} \end{aligned}$$

want  $P(\text{low}) \leq \delta$