

CS7800: Advanced Algorithms

Lecture 21 : Streaming II
- Finish Distinct Elements

Jonathan Ullman

12-2-2022

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 $IP(A_T = \text{low}) \geq 1 - \delta$

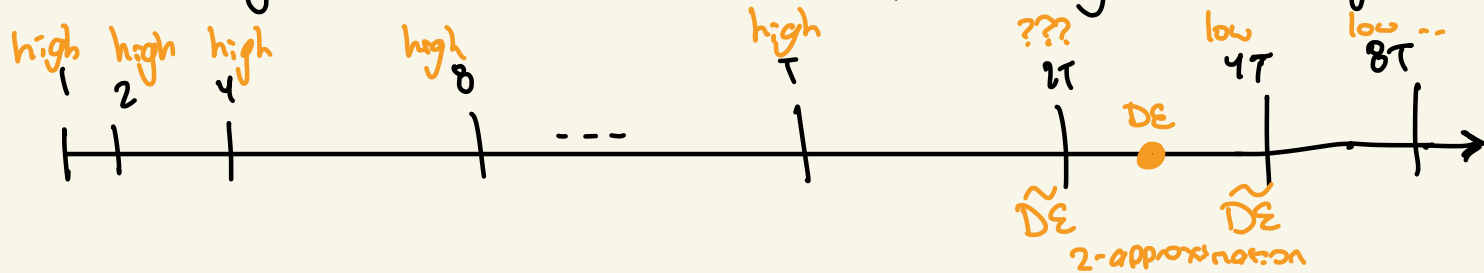
② If $DE \geq 2T$ then $IP(A_T = \text{high}) \geq 1 - \delta$

2 is arbitrary $\rightarrow 1 + \epsilon$

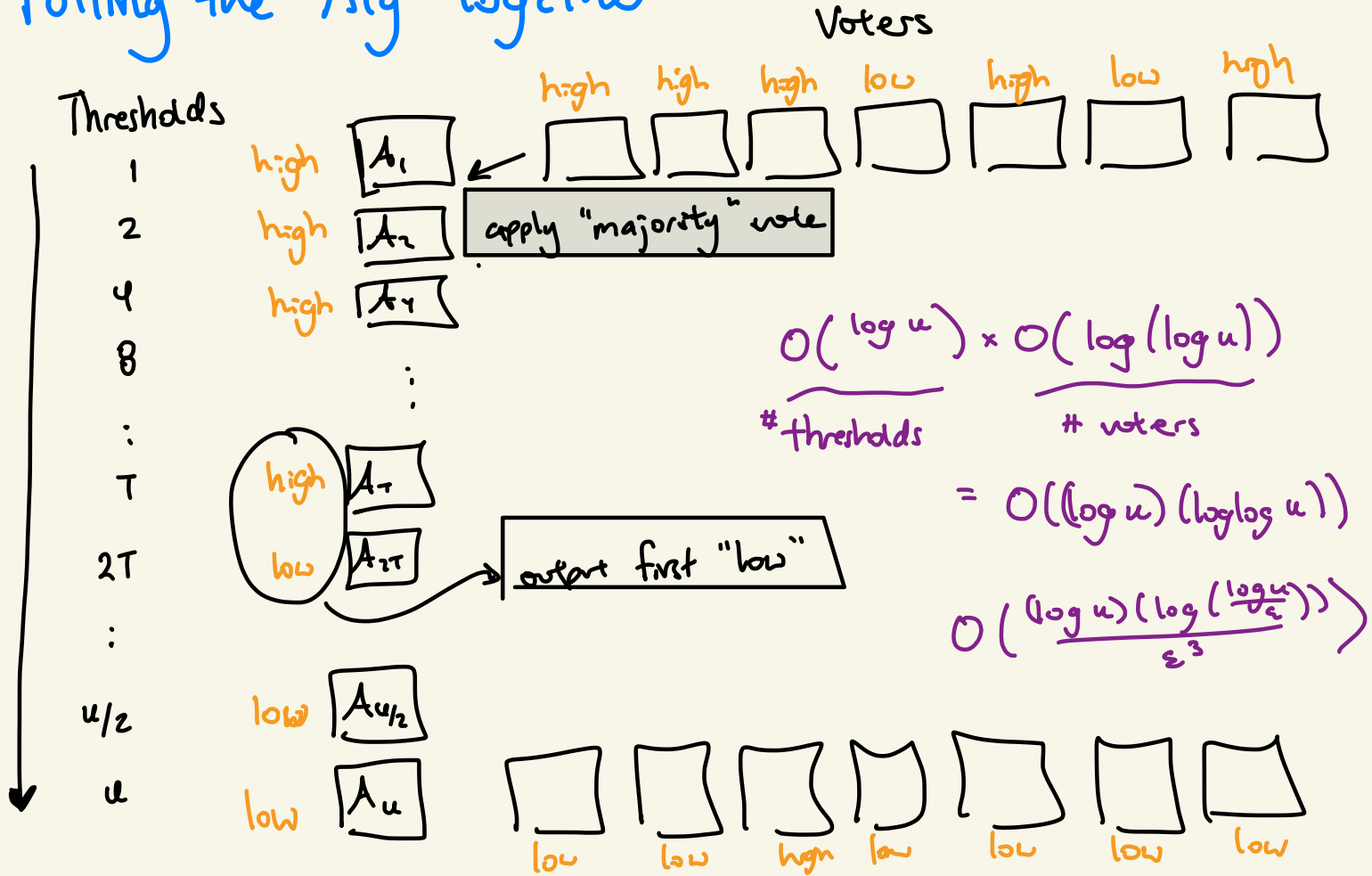
$$\log_{1+\epsilon} |U| = \frac{\log_2 |U|}{\log_2(1+\epsilon)} \approx \frac{\log_2 |U|}{\epsilon}$$

Can run $\log_2 |U|$ copies of $A_1, A_2, A_4, A_8, \dots, A_{|U|}$ in

parallel to get a 2-approximation with probability $\geq 1 - \delta \cdot \log_2 |U|$



Putting the Alg Together



Threshold Testing Distinct Elements I

Choose a ^{uniformly} 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 $DE \leq T$

$$P(\text{low}) = \left(1 - \frac{1}{T}\right)^{DE} \approx \frac{1}{e} \approx .36$$

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

Suppose $DE \geq 2T$

$$\begin{aligned} P(\text{low}) &= \left(1 - \frac{1}{T}\right)^{DE} \\ &= \left(\left(1 - \frac{1}{T}\right)^T\right)^2 \approx \frac{1}{e^2} \approx .14 \end{aligned}$$

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

Amplifying Success

$$\frac{1}{e} \approx .36 \quad \frac{1}{e^2} \approx .14 \quad \frac{1}{2}\left(\frac{1}{e} + \frac{1}{e^2}\right) \approx .25$$

Chernoff Bound: If z_1, \dots, z_n are independent $\{0,1\}$ random variables,

$z = z_1 + \dots + z_n$ and $\mu = \mathbb{E}(z)$ then

$$\mathbb{P}(z > \mu + \delta n) \leq e^{-n\delta^2/4}$$

$$\mathbb{P}(z < \mu - \delta n) \leq e^{-n\delta^2/4}$$

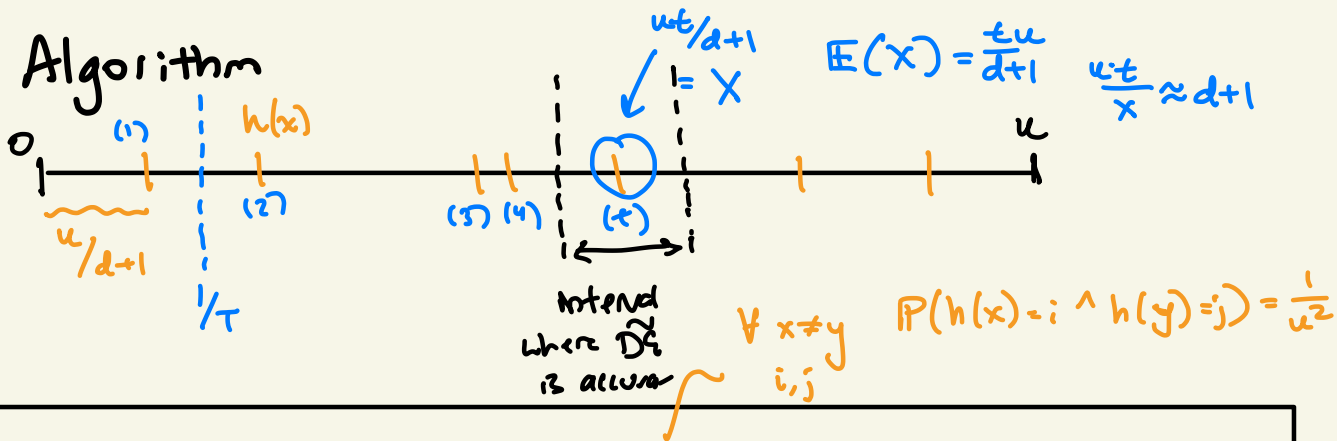
(low Dε)

$$\mu = .36n$$

$$\begin{aligned} & \mathbb{P}(z < .25n) \\ = & \mathbb{P}(z < \mu - .11n) \leq e^{-n(.11)^2/4} \approx e^{-\frac{n}{400}} \quad (\text{want}) \\ & \leq \beta \end{aligned}$$

$$n \geq 400 \cdot \ln(1/\beta) = O(\ln(1/\beta))$$

A Better Algorithm



- Let $\mathcal{H} = \{ h: [u] \rightarrow [u] \}$ be pairwise independent
- Choose h from \mathcal{H} randomly
- Let $t = \frac{1000}{\epsilon^2}$ and store the t smallest distinct hashes $(x_i, h(x_i))$
→ If you don't see t then just count exactly
- Let X be the t^{th} smallest $h(x_i)$ you've seen
- Return $\tilde{D}_\epsilon = \frac{t \cdot u}{X} (\approx d+1)$

Analysis

Thm: $P(|\tilde{D}\epsilon - D\epsilon| > \epsilon D\epsilon) \leq \frac{1}{100}$

- Let $\mathcal{H} \subseteq \{h: [u] \rightarrow [w]\}$ be pairwise independent
- Choose h from \mathcal{H} randomly
- Let $t = \frac{1000}{\epsilon^2}$ and store the t smallest distinct hashes $(x_i, h(x_i))$
- Let X be the t^{th} smallest $h(x_i)$ you've seen
- Return $\tilde{D}\epsilon = [\dots]$