

23 Feb 2026

## Distinct Elements Problem

Announcements:

- ① Weds 2/25: Eshan C guest lecturer
- ② Weds quiz covers 2/11, 2/18, 2/23

## Streaming model of computation

An algorithm processes a sequence of  $m$  tokens

from a universe of

$N$  potential tokens

$$N = 2^b$$

tokens  $\in \{0,1\}^b$

using storage space  $O(\text{poly}(b + \log m))$

and must afterward report an estimate of some function of the token sequence.

- Ex.
- most frequently occurring token
  - # distinct tokens
  - entropy of token distrib
  - ... et cetera

An  $(\epsilon, \delta)$ -PAC algorithm satisfies the property

$$\Pr(|\text{ALG's answer} - \text{TRUTH}| > \epsilon \cdot \text{TRUTH}) \leq \delta.$$

## Estimating distinct elements

Suppose we have stream

$$\theta_1, \theta_2, \dots, \theta_m \in \{0, 1\}^b$$

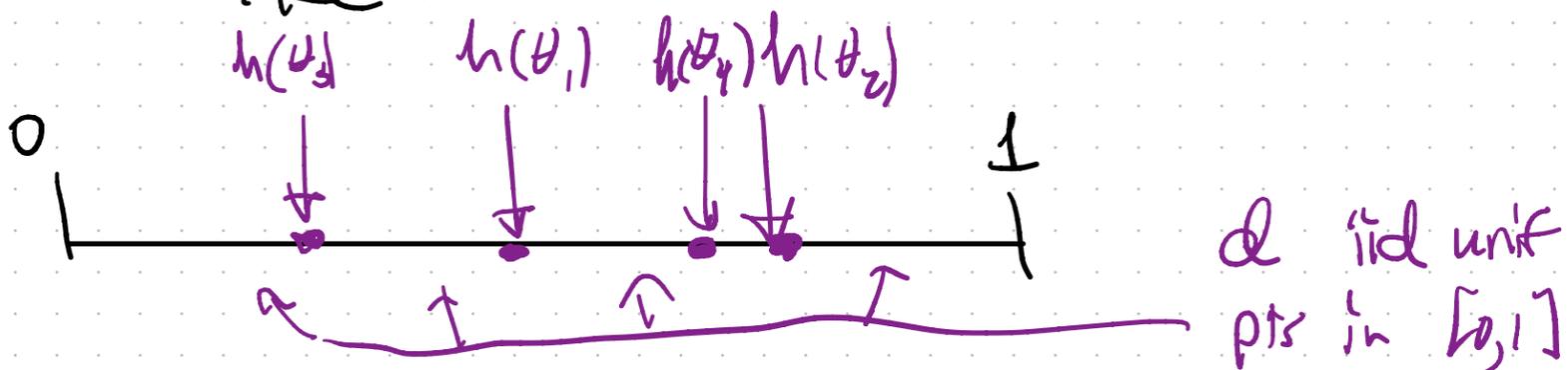
and (unrealistic) we have mutually independent random

$$h(\theta_1), h(\theta_2), \dots, h(\theta_m) \in [0, 1]$$

$$\text{if } \theta_i = \theta_j \quad h(\theta_i) = h(\theta_j)$$

$$\text{else } h(\theta_i) \text{ indep of } h(\theta_j).$$

If there are  $d$  distinct tokens in the stream, what would the set of distinct hash vals look like?



For  $d$  iid unif pts in  $[0,1]$

$$E[\min] = \frac{1}{d+1}.$$

## Distinct Elements Estimator (part 1)

Initialize 2-uniform hash  $h: [N] \rightarrow [p]$

( $p$  large prime)  $\leftarrow O(\log p)$  bits of storage

for each token  $\theta_i$ , compute  $h(\theta_i)$

and store  $Z = \min \{h(\theta_1), \dots, h(\theta_m)\}.$

// memory requirement is only  
storing  $Z$  and the  
most recent hash  $h(\theta_i)$ .

// If  $d$  distinct tokens,

expect  $Z \approx \frac{p}{d+1}.$

After seeing entire stream,  
report  $\hat{d} = \frac{p}{Z}.$

Claim. If  $p \gg d$

$$\Pr(\hat{d} \notin [\frac{1}{6}d, 6d]) < \frac{1}{3}.$$

Proof. Two bad events

$$\mathcal{E}_1: \hat{d} < \frac{d}{6} \iff Z > \frac{6p}{d}$$

$$\mathcal{E}_2: \hat{d} > 6d \iff Z < \frac{p}{6d}$$

$\Pr(\mathcal{E}_2)$ : the event  $Z < \frac{p}{6d}$  happens

when at least one of

$$\{h(\theta_1), \dots, h(\theta_m)\} \in \{1, \dots, \lfloor \frac{p}{6d} \rfloor - 1\}$$

Assume wlog  $\theta_1, \dots, \theta_d$  are the  $d$  distinct stream elements.

$$\Pr(\{h(\theta_1), \dots, h(\theta_m)\} \text{ intersects } \{1, \dots, \lfloor \frac{p}{6d} \rfloor - 1\})$$

$$\leq \mathbb{E} \left[ \# i \in \{1, \dots, d\} : h(\theta_i) \in \{1, \dots, \lfloor \frac{p}{6d} \rfloor - 1\} \right]$$

$$= \sum_{i=1}^d \Pr(h(\theta_i) < \frac{p}{6d})$$

$$< d \cdot \frac{1}{6d} = \frac{1}{6}.$$

$$\Pr(\mathcal{E}_2) < \frac{1}{6}$$

$\Pr(\mathcal{E}_1)$ :  $\mathcal{E}_1$  happens when

$\{h(\theta_1), h(\theta_2), \dots, h(\theta_d)\}$  disjoint

from  $\{1, 2, \dots, \frac{6p}{d}\}$ .

Step 1.  $\Pr(h(\theta_1) \notin \{1, \dots, \frac{6p}{d}\}) = \frac{p - 6p/d}{p}$

$$= 1 - \frac{6}{d}$$

$$\Pr(h(\theta_2) \notin \{1, \dots, \frac{6p}{d}\}) = 1 - \frac{6}{d}$$

⋮

$$\Pr(h(\theta_d) \notin \{1, \dots, \frac{6p}{d}\}) = 1 - \frac{6}{d}$$

Assumes  $h(\theta_1), \dots, h(\theta_d)$  mutually indep.

$$\Pr(\{h(\theta_1), \dots, h(\theta_d)\} \text{ disj from } \{1, \dots, \frac{6p}{d}\}) = \left(1 - \frac{6}{d}\right)^d$$

$$< \left[e^{-6/d}\right]^d = e^{-6}.$$

Lem. When  $X_1, \dots, X_n$  are  
pairwise indep't,

$$\text{Var}(X_1 + \dots + X_n) = \sum_{i=1}^n \text{Var}(X_i).$$

Proof. Let  $\bar{X}_i = \mathbb{E}(X_i)$ .

$$\begin{aligned}\bar{X} &= \bar{X}_1 + \dots + \bar{X}_n \\ &= \mathbb{E}[X_1 + \dots + X_n].\end{aligned}$$

$$\begin{aligned}\text{Var}(X_1 + \dots + X_n) &= \mathbb{E}\left[\left((X_1 + \dots + X_n) - \bar{X}\right)^2\right] \\ &= \mathbb{E}\left[\sum_{i=1}^n X_i^2 + 2 \sum_{1 \leq i < j \leq n} X_i X_j \right. \\ &\quad \left. - 2 \sum X_i \bar{X} + \bar{X}^2\right]\end{aligned}$$

$$\sum \text{Var}(X_i) = \sum_{i=1}^n \mathbb{E} \left[ X_i^2 - 2X_i\bar{X} + \bar{X}^2 \right]$$

Abandoning proof to avoid  
careless algebra error.

Theme of proof.  $\text{Var}(X_1 + \dots + X_n)$

is  $\mathbb{E}[\text{quadratic polynomial}]$ .

Each monomial depends

on  $\leq 2$  variables, so

can't detect whether

they are fully indep.

or just pairwise indep.

Applying to

$$Pr(\mathcal{E}_1) = Pr(\{h(\theta_1), \dots, h(\theta_d)\} \text{ disj})$$

from  $\{1, \dots, \frac{b}{d}\}$

$$i=1, \dots, d:$$
$$X_i = \begin{cases} 1 & \text{if } h(\theta_i) \in \{1, \dots, \frac{b}{d}\} \\ 0 & \text{o.w.} \end{cases}$$

$$E(X_i) = \frac{b}{d}$$

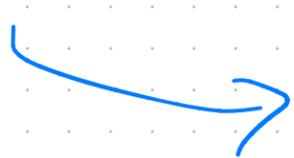
$$\text{Var}(X_i) = \frac{b}{d} \left(1 - \frac{b}{d}\right) < E[X_i]$$

$$\mathcal{E}_1 = \{X_1 + \dots + X_d = 0\}$$

Have:  $E(X_1 + \dots + X_d) = \frac{6}{d} \cdot d = 6$

$$\text{Var}(X_1 + \dots + X_d)$$

using  
pairwise  
indep.



$$= \text{Var}(X_1) + \dots + \text{Var}(X_d)$$

$$< E(X_1) + \dots + E(X_d)$$

$$= 6$$

$$Pr(X_1 + \dots + X_d = 0)$$

Chebyshev



$$\frac{\text{Var}(X_1 + \dots + X_d)}{6^2}$$



$$\frac{6}{6^2}$$

$$= \frac{1}{6}$$

$$Pr(\mathcal{E}_1), Pr(\mathcal{E}_2) \leq \frac{1}{6}$$

$$\Pr(\varepsilon_1 \cup \varepsilon_2) \leq \frac{1}{3}.$$

$$\Pr(\hat{d} \in \left[\frac{1}{6}d, 6d\right]) \geq \frac{2}{3}.$$

Reducing relative error to  $\pm \varepsilon$ .

Let  $z_1, z_2, \dots, z_t$  be the  $t$  smallest hash values.

$$t = \left\lceil \frac{2}{\varepsilon^2 \delta} \right\rceil.$$

$$\hat{d} = \frac{t \cdot p}{z_t} \quad \begin{array}{l} z_1 \approx \frac{p}{d+1} \\ z_t \approx \frac{tp}{d+1} \end{array}$$

Analysis in § 3.2