

22 Apr 2022

Streaming Algorithm for Distinct Elements

Input stream a_1, \dots, a_n .

Each a_i belongs to $\{0, 1\}^b$ (set of "tokens")

$m = 2^b = \#$ of possible tokens.

Algorithm with working space $S = \text{poly}(\log n, \log m)$

observes a_1, \dots, a_n one by one, and afterward must estimate $\#$ distinct tokens appearing in the sequence.

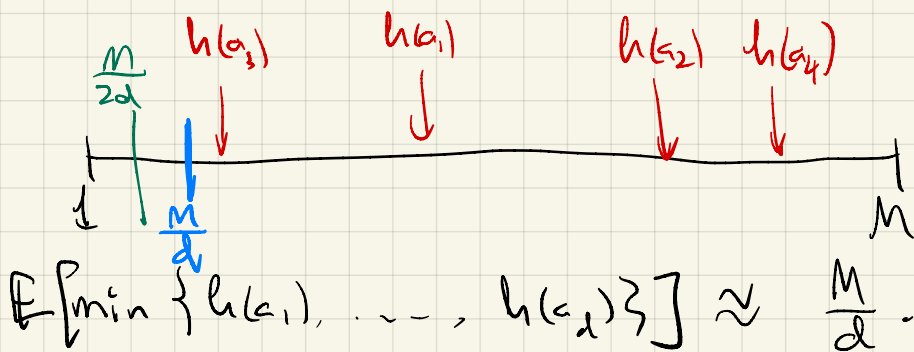
GOAL: (ϵ, δ) -PAC meaning: if true number of distinct tokens is d , then

$$\Pr(\text{ALG's output is not in } [(1-\epsilon)d, (1+\epsilon)d]) < \delta.$$

Use hash function $h: \{0, 1\}^b \rightarrow [M]$

for some large integer M . (Typically $M \geq 2^b$.)

Suppose for now that a_1, \dots, a_d are the d distinct tokens that appear in the stream and that $h(a_1), h(a_2), \dots, h(a_d)$ are independent uniformly random elements of $[M]$.



$$\Pr(h(a_1), \dots, h(a_d) > \frac{M}{d}) = \left(1 - \frac{1}{d}\right)^d \approx e^{-1}$$

$$\begin{aligned} \Pr(\min\{h(a_i)\} \leq \frac{M}{2d}) &\leq \sum_{i=1}^d \Pr(h(a_i) \leq \frac{M}{2d}) \\ &= \frac{1}{2}. \end{aligned}$$

- Alg.
- ① Initialize $z = M$
 - ② for each a_i in succession
 - compute $h(a_i)$
 - if $h(a_i) < z$, $z \leftarrow h(a_i)$
 - // invariant: after t iterations
 $z = \min \{ h(a_j) \mid 1 \leq j \leq t \}$
 - ③ output M/z as our estimate of d .

Analysis. For $k \in [M]$ let $X_{ik} = \begin{cases} 1 & \text{if } h(a_i) \leq k \\ 0 & \text{if } h(a_i) > k \end{cases}$

Fact 1. $E[X_{ik}] = \frac{k}{M}$

$Y_k = \sum_{i=1}^d X_{ik}$
 = # distinct tokens in stream whose hash value is $\leq k$.
 assuming a_1, \dots, a_d are the distinct tokens

Fact 2. $E[Y_k] = \frac{dk}{M}$

① $h(a_i)$ is uniform in $[M]$

Fact 3. $\text{Var}(Y_k) < \frac{dk}{M}$ for all k

②

$h(a_i)$ and $h(a_j)$ are independent

$$\begin{aligned} \text{Var}(Y_k) &= E(Y_k^2) - E(Y_k)^2 \\ &= \sum_{i=1}^d \sum_{j=1}^d E[X_{ik} X_{jk}] - \left(\frac{dk}{M}\right)^2 \\ &= \sum_{i=1}^d E[X_{ik}^2] + 2 \sum_{1 \leq i < j \leq d} E[X_{ik} X_{jk}] - \left(\frac{dk}{M}\right)^2 \end{aligned}$$

$$\begin{aligned}
&= \frac{dk}{M} + 2 \binom{d}{2} \left(\frac{k}{M}\right)^2 - \left(\frac{dk}{M}\right)^2 \\
&= \frac{dk}{M} + (d^2 - d) \left(\frac{k}{M}\right)^2 - \left(\frac{dk}{M}\right)^2 \\
&= \frac{dk}{M} - \frac{dk^2}{M^2} \\
&< \frac{dk}{M}.
\end{aligned}$$

Def. A 2-universal hash function (UHF) is a probability distribution over hash functions that satisfies:

- I. $h(a_i)$ is unif. distrib. over its range for all i .
- II. $h(a_i), h(a_j)$ independent $\forall i \neq j$.
"pairwise independence"

Example. Say M is prime, and identify $\{0, 1\}^b$ with $[m]$, and assume $M \geq m$.

Define h as follows: $h_{ab}(x) = ax + b \pmod{M}$
 Sample a, b as indep't, unif random mod M , set $h = h_{ab}$.

Storing a description of h takes
 $2 \log(M)$ bits to store a, b .

To verify 2-universality, must show

$$\forall a_i \neq a_j \quad \Pr_{a,b} \left(h(a_i) = x, h(a_j) = y \right) = \frac{1}{M^2} \quad \forall x, y$$

$$h_{ab}(a_i) = x, \quad h_{ab}(a_j) = y \quad \text{means}$$

$$\begin{bmatrix} a_i & 1 \\ a_j & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} \equiv \begin{bmatrix} x \\ y \end{bmatrix} \pmod{M}$$

The linear system has a unique solution

\pmod{M} b/c $\mathbb{Z}/(M)$ is a field

$$\text{and} \quad \det \begin{pmatrix} a_i & 1 \\ a_j & 1 \end{pmatrix} = a_i - a_j \not\equiv 0 \pmod{M}$$