

4.5 Streams and Sketches

Thursday, January 30, 2020 12:01 AM

Consider a sequence $a_1, \dots, a_n \in [m]$, with n and m very large.

We want to compute $|\text{uniq}\{a_1, \dots, a_n\}|$.

$O(m)$ solution - bit vector

$O(n \log m)$ solution - store list of items seen

(don't need to organize if we only care about time complexity)

Goal: $O(\log n \cdot \log m)$.

Thm: (lower bound) Any exact deterministic algorithm must use at least m bits of memory on some sequence of length $m+1$.

Proof. Suppose we have seen a_1, \dots, a_m already.

Now assume for contradiction that our algorithm uses $\leq m$ bits of memory on all such sequences.

$|\mathcal{P}([m])| = 2^m$, and $\text{uniq}\{a_1, \dots, a_m\}$ can be any subset except \emptyset ,

so $2^m - 1$ possibilities for $\text{uniq}\{a_1, \dots, a_m\}$.

But we only have 2^{m-1} possible memory states.

So two different subsets $S_1, S_2 \in \mathcal{P}([m]) \setminus \emptyset$ must have the same memory state.

$|S_1| = |S_2|$ because otherwise our algorithm must be wrong for one of them.

Now add an $a_{m+1} \in S_1$. Then the new sets would have the same state but different sizes.



Consider an idealized streaming algorithm (ISA)

1. Pick random hash function $h: [m] \rightarrow [0, 1]$.

2. Calculate $z = \min_{i \in \text{stream}} h(i)$.

3. Output $\frac{1}{z} - 1$

Let $\mathcal{S} = \text{unique}\{a_1, \dots, a_n\} = \{j_1, \dots, j_t\}$

$h(j_1), \dots, h(j_t) = X_1, \dots, X_n$ be ind. Unif $[0, 1]$

$z = \min \{X_i\}_{i=1}^n$

Lemma: If $X: \Omega \rightarrow [0, +\infty]$ is a nonnegative r.v., then

$$E X = \int_{[0, +\infty)} P(X > x) dx$$

proof. For every $\omega \in \Omega$, $X(\omega) = \int_{(0, X(\omega))} dx = \int_{[0, +\infty)} \mathbb{1}_{(0, X(\omega))}(x) dx$

$$E X = \int_{\Omega} X dP = \int_{\Omega} \int_{[0, +\infty)} \mathbb{1}_{(0, X(\omega))}(x) dx dP(\omega)$$

$$= \int_{[0, +\infty)} \int_{\Omega} \mathbb{1}_{(0, X(\omega))}(x) dP(\omega) dx = \int_{[0, +\infty)} P(X > x) dx$$

If $x < X(\omega)$, $\mathbb{1} = 1$.

Integral of $1 \cdot P(\omega)$ for all ω st. $x < X(\omega)$.



Claim: $E Z = \frac{1}{t+1}$

proof. $E Z = \int_0^{\infty} P(Z > \lambda) d\lambda = \int_0^1 P(\forall i, X_i > \lambda) d\lambda$
 $= \int_0^1 \prod_{i=1}^t P(X_i > \lambda) d\lambda = \int_0^1 (1-\lambda)^t d\lambda = \left[\frac{(1-\lambda)^{t+1}}{t+1} \right]_0^1 = \frac{1}{t+1}$

Claim: $E Z^2 = \frac{2}{(t+1)(t+2)}$

proof. $E Z^2 = \int_0^1 P(Z^2 > \lambda) d\lambda = \int_0^1 P(Z > \sqrt{\lambda}) d\lambda$
 $= \int_0^1 (1-\sqrt{\lambda})^t d\lambda = 2 \int_1^0 u^t (u-1) du = 2 \int_0^1 u^t (1-u) du$

Let $u = 1 - \sqrt{\lambda}$, $\sqrt{\lambda} = 1 - u$

$du = -\frac{1}{2\sqrt{\lambda}} d\lambda \Rightarrow 2 du = \frac{1}{u-1} d\lambda$

$d\lambda = 2(u-1) du$

$\bar{u} = (1-u) \quad dv = u^t du$

$d\bar{u} = -du \quad v = \frac{1}{t+1} u^{t+1}$

$= 2 \left[(1-u) \cdot \frac{1}{t+1} u^{t+1} \Big|_0^1 - \int_0^1 \frac{1}{t+1} u^{t+1} du \right]$

$= \frac{2}{(t+1)(t+2)}$



Thus, $\text{Var}[Z] = E Z^2 - (E Z)^2 = \frac{2}{(t+1)(t+2)} - \frac{1}{(t+1)^2} < \frac{1}{(t+1)^2}$

$$\text{Thus, } \text{Var}[Z] = \mathbb{E} Z^2 - (\mathbb{E} Z)^2 = \frac{t}{(t+1)^2(t+2)} < \frac{1}{(t+1)^2}$$

Averaging algorithm

1. Run $q = \frac{1}{\epsilon^2 \eta}$ ISAs in parallel

$$2. \bar{z} = \frac{1}{q} \sum_{i=1}^q z_i$$

3. output $\frac{1}{\bar{z}} - 1$

$$\text{Then } \mathbb{E}(\bar{z}) = \frac{1}{t+1}, \quad \text{Var}(\bar{z}) = \frac{1}{q} \cdot \frac{t}{(t+1)^2(t+2)} < \frac{1}{q(t+1)^2}.$$

$$\text{By Chebyshev } \mathbb{P}\left(\left|\bar{z} - \frac{1}{t+1}\right| > \frac{\epsilon}{t+1}\right) < \frac{(t+1)^2}{\epsilon^2} \cdot \frac{1}{q(t+1)^2} = \eta.$$

$$\text{Claim: } \mathbb{P}\left(\left|\left(\frac{1}{\bar{z}} - 1\right) - t\right| > O(\epsilon)t\right) < O(\eta)$$

$$\mathbb{P}\left(\left|\bar{z} - \frac{1}{t+1}\right| > \frac{\epsilon}{t+1}\right) < \eta$$

$$\mathbb{P}\left(\left|\bar{z}t + \bar{z} - 1\right| > \epsilon\right) < \eta$$

$$\mathbb{P}\left(\left|\frac{1}{\bar{z}} - t - 1\right| > \frac{\epsilon}{|\bar{z}|}\right) < \eta$$

$$\text{w.p. } 1-\eta \quad |\bar{z}| \leq \frac{1+\epsilon}{t+1},$$

$$\Rightarrow \mathbb{P}\left(\left|\frac{1}{\bar{z}} - t - 1\right| > \underbrace{\frac{\epsilon(t+1)}{1+\epsilon}}_{O(\epsilon)t} \right) < \underbrace{\eta}_{2\eta}.$$

Chebyshev

$$\mathbb{P}\left(|x - \mathbb{E}x| \geq a\right) \leq \frac{\text{Var}(x)}{a^2}$$

$$\epsilon(t+1)\left(1 - \epsilon + \frac{\epsilon^2}{2} - O(\epsilon^3)\right)$$

$$= \epsilon + t\epsilon + O(\epsilon^2)$$

$$= O(\epsilon)t.$$



So, with high prob., our estimator is within a factor $(1+O(\epsilon))$ of t .

But we have forgotten to take into account the $\geq n$ bits needed to store our hash function.

Let \mathcal{H} be a set of functions that map $[a] \rightarrow [b]$.

Define: \mathcal{H} is a k -wise independent hash family if

$$\forall i_1 \neq i_2 \neq \dots \neq i_k \in [a]$$

$$\text{and } \forall j_1, \dots, j_k \in [b],$$

$$\mathbb{P}_{h \sim \mathcal{H}} (h(i_1) = j_1 \wedge \dots \wedge h(i_k) = j_k) = \frac{1}{b^k}.$$

Ex. The set \mathcal{H} of all functions $[a] \rightarrow [b]$ is k -wise indep for every k .
 $|\mathcal{H}| = b^a$, so $h \in \mathcal{H}$ is representable in a $\lg b$ bits.

Ex. Let $a = b = q = \text{prime power}$. \mathcal{H} will be the set of all deg $k-1$ polynomials in $\mathbb{F}_q[x]$.

Claim: $\mathcal{H}_{\text{poly-}k}$ is a k -wise family.


proof. Lagrange interpolation. If we know i_1, \dots, i_k and j_1, \dots, j_k and that no i 's repeat, then

$$p(x) = \sum_{r=1}^k \left(\frac{\prod_{y \in [k] \setminus \{r\}} (x - i_y)}{\prod_{y \in [k] \setminus \{r\}} (i_r - i_y)} \right) \cdot j_r$$

satisfies $\forall r \ p(i_r) = j_r$ and this polynomial is unique because \mathbb{F}_q is a field.

$$\text{Thus, } |\mathcal{H}_{\text{poly-}k}| = q^k.$$

But note that there is precisely one poly of deg $\leq k-1$ s.t.

it goes through all (i_r, j_r) , so $\mathbb{P}_{h \sim \mathcal{H}_{\text{poly-}k}} (h(i_1) = j_1 \wedge \dots \wedge h(i_k) = j_k) = \frac{1}{q^k}$. 

Also, each $h \in \mathcal{H}_{\text{poly-}k}$ is representable using $k \log q$ bits.

How much independence do we need for unique items?

How much independence do we need for unique items?

We can replace independent hash functions with pairwise independence

Careful analysis in book, but basically pairwise independence allows us to sum variances.