

11. Count distinct sketch

Wednesday, October 6, 2021 2:33 PM

Last time: How to construct random hash families

Today: Frequency moments + estimating F_0 .

Consider a sequence $a_1, \dots, a_n \in [m]$, with n, m both large (stream)

$\forall s \in [m]$, let $f_s = |\{i \mid a_i = s\}|$, the frequency of s in stream.

Def. For $p \in \mathbb{N}$, the p th freq moment of the stream is

$$F_p = \sum_{s=1}^m f_s^p.$$

($0^0 = 0$)

Aside: $F_0 = |\text{Uniq}\{a_1, \dots, a_n\}|$, number of distinct symbols

Aside: $F_1 = n$

Aside: F_2 is useful in computing variance in occurrences of elements.

$$\begin{aligned} \frac{1}{m} \sum_{s=1}^m (f_s - \frac{n}{m})^2 &= \frac{1}{m} \sum_{s=1}^m (f_s^2 - 2 \frac{n}{m} f_s + (\frac{n}{m})^2) = \left(\frac{1}{m} \sum_{s=1}^m f_s^2 \right) - \frac{2n}{m^2} \sum_{s=1}^m f_s + \frac{n^2}{m^2} \\ &= \frac{F_2}{m} - \frac{n^2}{m^2}. \end{aligned}$$

Aside: $\lim_{p \rightarrow \infty} F_p^{1/p} = \lim_{p \rightarrow \infty} \left(\sum_{s=1}^m f_s^p \right)^{1/p}$ is the freq of the most frequent element(s).

F_0 -sketching (distinct elements)

$O(m)$ space solution - bit vector

$O(n \log m)$ space sol - store a list of items seen

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

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

proof. Assume that an alg uses $< m$ bits of memory on all such seq.

Recall that the power set $|\mathcal{P}([m])| = 2^m$, and $\text{Uniq}\{a_1, \dots, a_m\}$

Recall that the power set $|\mathcal{P}([m])| = 2^m$, and $\text{uniq}\{a_1, \dots, a_m\}$ can be any subset except \emptyset , so $2^m - 1$ possibilities.

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

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

$\Rightarrow |S_1| = |S_2|$ because otherwise our alg would be wrong.

Let $b \in S_1$ but $b \notin S_2$. Then $S_1 \cup \{b\} = S_1$ so $|S_1 \cup \{b\}| = |S_1|$

and $S_2 \cup \{b\}$, so $|S_2 \cup \{b\}| = |S_2| + 1$.

But the alg will have same memory state after adding b , so it must be wrong on one of them. □

Consider an idealized streaming alg (ISA) [Telai Nelson's book on sketching]

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

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

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

Let $\mathcal{A} = \text{uniq}\{a_1, \dots, a_n\} = \{b_1, \dots, b_d\}$

$h(b_1), \dots, h(b_d) = X_1, \dots, X_d$ i.i.d. $\text{Unif.}[0, 1]$

$$z = \min_{i=1}^d \{X_i\}$$

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

$$\mathbb{E}X = \int_{[0, +\infty)} \text{Prob}(X > x) dx.$$

Claim: $\mathbb{E}z = \frac{1}{d+1}$

proof: $\mathbb{E}z = \int_0^{\infty} \text{Prob}(z > \lambda) d\lambda = \int_0^1 \text{Prob}(\forall i, X_i > \lambda) d\lambda$
 $= \int_0^1 \prod_{i=1}^d \text{Prob}(X_i > \lambda) d\lambda = \int_0^1 (1-\lambda)^d d\lambda = \left[-\frac{(1-\lambda)^{d+1}}{d+1} \right]_0^1 = \frac{1}{d+1}$ □

Claim: $\mathbb{E}z^2 = \frac{2}{(d+1)(d+2)}$

Claim: $\mathbb{E} z^2 = \frac{2}{(d+1)(d+2)}$

Proof:
$$\begin{aligned} \mathbb{E} z^2 &= \int_0^1 \text{Prob}(z^2 > \lambda) d\lambda = \int_0^1 \text{Prob}(z > \sqrt{\lambda}) d\lambda \\ &= \int_0^1 (1 - \sqrt{\lambda})^d d\lambda = 2 \int_0^1 u^d (u-1) du = 2 \int_0^1 u^d (1-u) du \\ &\quad \text{Let } u = 1 - \sqrt{\lambda} \curvearrowright \\ &= \frac{2}{(d+1)(d+2)} \end{aligned}$$



Then, $\text{Var}(z) = \mathbb{E} z^2 - (\mathbb{E} z)^2 = \frac{d}{(d+1)^2(d+2)} < \frac{1}{(d+1)^2}$.

Averaging algorithm (AA)

1. Run $q = \frac{1}{\epsilon^2 n}$ ISAs in parallel
2. $\bar{z} = \frac{1}{q} \sum_{i=1}^q z_i$, where z_i comes from ISAs
3. Output $\frac{1}{\bar{z}} - 1$.

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

By Chebyshev, $\text{Prob}\left(\left|\bar{z} - \frac{1}{d+1}\right| > \frac{\epsilon}{d+1}\right) < \frac{(d+1)^2}{\epsilon^2} \cdot \frac{1}{q(d+1)^2} = \eta$

Claim: $\text{Prob}\left(\left|\left(\frac{1}{\bar{z}} - 1\right) - d\right| > O(\epsilon)d\right) < O(\eta)$.

Proof:

$$\begin{aligned} \text{Prob}\left(\left|\bar{z} - \frac{1}{d+1}\right| > \frac{\epsilon}{d+1}\right) &< \eta \\ \text{Prob}\left(\left|\bar{z}d + \bar{z} - 1\right| > \epsilon\right) &< \eta \\ \text{Prob}\left(\left|\frac{1}{\bar{z}} - d - 1\right| > \frac{\epsilon}{|\bar{z}|}\right) &< \eta \end{aligned}$$

$$\begin{aligned} \epsilon(d+1)(1-\epsilon + \frac{\epsilon^2}{2} - O(\epsilon^3)) \\ = \epsilon + d\epsilon + O(\epsilon^2) \\ = O(\epsilon)d \end{aligned}$$

w.p. $1 - \eta$, $\left|\bar{z}\right| \leq \frac{1+\epsilon}{d+1}$

$$\Rightarrow \text{Prob}\left(\left|\frac{1}{\bar{z}} - d - 1\right| > \frac{\epsilon(d+1)}{1+\epsilon}\right) < \eta$$

$$\Rightarrow \text{Prob} \left(\left| \frac{1}{\bar{z}} - d - 1 \right| > \underbrace{\frac{\epsilon(d+1)}{1+\epsilon}}_{O(\epsilon)d} \right) < \underbrace{\eta}_{O(\eta)}$$



With high prob our estimator \bar{D} within a factor $1+O(\epsilon)$ of $1+d$.

Space-complexity $O\left(\frac{1}{\epsilon^2 \eta}, \text{space } [0,1]\right)$

Median of averages

1. Instantiate $s = \lceil 36 \ln\left(\frac{2}{\delta}\right) \rceil$ ind. copies of averaging alg with $\eta = \frac{1}{3}$
2. Output the median \hat{d} of $\left\{ \frac{1}{\bar{z}_j} - 1 \right\}_{j=1}^s$ where \bar{z}_j is the j th output of AA.

Claim: $\text{Prob} \left(\left| \hat{d} - d \right| > \epsilon d \right) < \delta$

Proof: Let $Y_j = \begin{cases} 1 & \text{if } \left| \frac{1}{\bar{z}_j} - d \right| > \epsilon d \\ 0 & \text{otherwise} \end{cases}$

The median fails if at least half of the $Y_j = 1$, i.e. $\sum_{j=1}^s Y_j > \frac{s}{2}$.

Note $\text{Prob} \left(\sum Y_j > \frac{s}{2} \right) = \text{Prob} \left(\sum Y_j - \frac{s}{3} > \frac{s}{6} \right)$

Simplify by using the stronger assumption that $\mathbb{E} Y_j = \frac{1}{3}$ (i.e. higher failure rate)

$$= \text{Prob} \left(\sum Y_j - \mathbb{E} \sum Y_j > \frac{1}{2} \mathbb{E} \sum Y_j \right) < \exp \left(\frac{-(0.5)^2 \cdot s/3}{3} \right) < \delta$$

\uparrow Chernoff bound



Space-complexity $O\left(\frac{1}{\epsilon^2} \lg \frac{1}{\delta}, \text{space } [0,1]\right)$

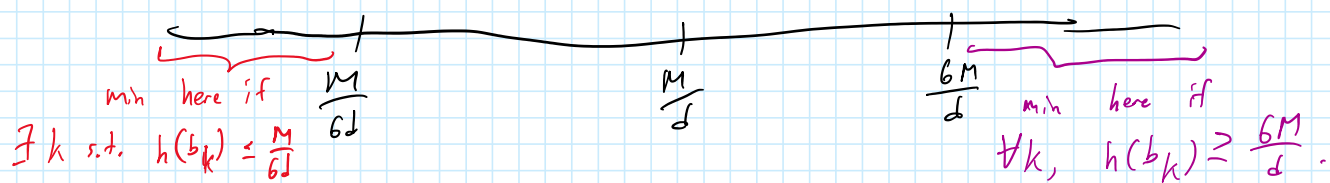
Consider the following algorithm [from book]

- Same notation as before
- Pick z -ind hash function $h: [m] \rightarrow [M]$, where $M > m$.
- Output $z = \frac{M}{\min_{i \in \text{stream}} h(a_i)}$, an estimate for d , # distinct items = min

Note: $S = \{h(b_1), \dots, h(b_d)\}$ is a set of d random and pairwise ind. uniform values from $[M]$.

Claim: With prob at least $\frac{2}{3} - \frac{d}{M}$, $\frac{d}{6} \leq \ell \leq 6d$.

proof: Equivalent to show $\frac{M}{6d} \leq \min \leq \frac{6M}{d}$.



$$\begin{aligned} \bullet \text{ Prob} \left(\min \leq \frac{M}{6d} \right) &= \text{Prob} \left(\exists k, h(b_k) \leq \frac{M}{6d} \right) \\ &\leq \sum_{i=1}^d \text{Prob} \left(h(b_i) \leq \frac{M}{6d} \right) \leq d \left(\frac{\frac{M}{6d}}{M} \right) \\ &\leq d \left(\frac{1}{6d} + \frac{1}{M} \right) \leq \frac{1}{6} + \frac{d}{M}. \end{aligned}$$

$$\bullet \text{ Prob} \left(\min \geq \frac{6M}{d} \right) = \text{Prob} \left(\forall k, h(b_k) \geq \frac{6M}{d} \right)$$

$$\text{let } y_i = \begin{cases} 0, & \text{if } h(b_i) \geq \frac{6M}{d} \\ 1, & \text{otherwise} \end{cases}, \quad y = \sum_{i=1}^d y_i.$$

$$= \text{Prob}(y = 0)$$

$$\begin{aligned} &\leq \text{Prob}(|y - \mathbb{E}y| \geq \mathbb{E}y) \\ \text{Chebyshev} \quad &\leq \frac{\text{Var}(y)}{(\mathbb{E}y)^2} \leq \frac{1}{\mathbb{E}y} \leq \frac{1}{6}. \end{aligned}$$

$$\text{Note: } \text{Prob}(y_i = 1) = \text{Prob} \left(h(b_i) < \frac{6M}{d} \right) \geq \frac{\frac{6M}{d}}{M} = \frac{6}{d}$$

$$\Rightarrow \mathbb{E}y_i \geq \frac{6}{d}, \quad \mathbb{E}y \geq 6$$

Also, $\text{Var}(y) = d \text{Var}(y_i)$ by 2-independence.

$$\text{Var}(y_i) = \mathbb{E}(y_i^2) - \mathbb{E}y_i^2 = \mathbb{E}y_i - \mathbb{E}y_i^2 \leq \mathbb{E}y_i$$

$$\Rightarrow \text{Var}(y) \leq \mathbb{E}y$$

Thus, with prob at least $\frac{2}{3} - \frac{d}{M}$, $\frac{d}{6} \leq \ell \leq 6d$.

