

12. More frequency moment sketching

Thursday, October 7, 2021 3:06 PM

Last time: Idealized F_0 sketch and constant approximation

Today: ϵ -approximation F_0 sketch and F_2 sketch

Non-idealized ^(Geometric sampling) Hajnal-Martin counting [1985]

1. Pick h from a 2-wise family $[n] \rightarrow [n]$ for n a power of 2.
2. Maintain $X = \max_{a_i \in \text{stream}} \text{lsb}(h(a_i))$, where lsb is the least-significant 1-bit of a number
3. Output $\tilde{d} = 2^X$

unique $\{a_1, \dots, a_n\}$

$$\begin{aligned} \text{lsb}(101101100) &= 3 \\ \text{lsb}(10110000) &= 5 \end{aligned}$$

For fixed j , let $Z_j = |\{b_i \in \text{stream} \mid \text{lsb}(h(b_i)) = j\}|$

Let $Z_{>j} = |\{b_i \in \text{stream} \mid \text{lsb}(h(b_i)) > j\}|$

$$\text{Let } Y_i = \begin{cases} 1 & \text{if } \text{lsb}(h(b_i)) = j \\ 0 & \text{otherwise} \end{cases} \quad \mathbb{E} Y_i = \frac{1}{2^{j+1}}$$

$$\text{Var}(Y_i) = \mathbb{E} Y_i^2 - (\mathbb{E} Y_i)^2 = \frac{1}{2^{j+1}} - \frac{1}{2^{2j+2}} < \frac{1}{2^{j+1}}$$

$$\text{Then } Z_j = \sum_{b_i \in \text{str}} Y_i. \quad \mathbb{E} Z_j = \frac{d}{2^{j+1}}, \text{ where } d = |\text{uniq}(\text{stream})|$$

$$\mathbb{E} Z_{>j} = d \left(\frac{1}{2^{j+2}} + \frac{1}{2^{j+3}} + \dots \right) = \frac{d}{2^{j+1}}$$

$$\begin{aligned} \text{Var}(Z_j) &= \text{Var}\left(\sum Y_i\right) = \mathbb{E}\left(\sum Y_i\right)^2 - \left(\mathbb{E}\sum Y_i\right)^2 \\ &= \mathbb{E}\left[\sum Y_i^2 + 2 \sum_{i \neq i_2} Y_i Y_{i_2}\right] = \sum \mathbb{E} Y_i^2 + 2 \sum_{i \neq i_2} \mathbb{E} Y_i Y_{i_2} \\ &= \sum \left[\mathbb{E} Y_i^2 - (\mathbb{E} Y_i)^2 \right] + 2 \sum_{i \neq i_2} (\mathbb{E} Y_i)(\mathbb{E} Y_{i_2}) \\ &= \sum \text{Var}(Y_i) + 2 \sum_{i \neq i_2} (\mathbb{E} Y_i)(\mathbb{E} Y_{i_2}) = \sum \text{Var}(Y_i) + 2 \sum_{i \neq i_2} \frac{1}{2^{2j+2}} \\ &< \frac{d}{2^{j+1}} + \frac{2d^2}{2^{2j+2}} \Rightarrow \text{Var}(Z_j) < \frac{d}{2^{j+1}} \end{aligned}$$

Now for $j^* = \lfloor \lg d - 5 \rfloor$, we have

$$\lg d - 6 \leq j^* \leq \lg d - 5$$

$$\text{Thus } \frac{d}{2^{5d-4}} \leq \mathbb{E} Z_{j^*} \leq \frac{d}{2^{\lg d - 5}}$$

$$\Rightarrow 16 \leq \mathbb{E} Z_{j^*} \leq 32$$

$$\frac{d}{2^{j^*}}$$

$$\Rightarrow 16 \leq \mathbb{E} z_{j^*} \leq 32$$

$$\text{Prob}(z_{j^*} = 0) \leq \text{Prob}(|z_{j^*} - \mathbb{E} z_{j^*}| \geq 16) \leq \frac{\text{Var}(z_{j^*})}{256} < \frac{2^{\frac{d}{\lg d - 5}}}{256} = \frac{1}{8}$$

↑
Chebyshev.

For $j = \lfloor \lg d + 5 \rfloor$, we have $j > \lg d + 4$

$$\mathbb{E} z_{>j} \leq \frac{d}{2^{\lg d + 5}} = \frac{1}{32} \Rightarrow \text{Prob}(z_{>j} \geq 1) < \frac{1}{32} \text{ by Markov}$$

Thus w.p. $1 - \frac{1}{8} - \frac{1}{32}$, the max lsb will be b/t j^* and j , in a constant range,

i.e. we get a $O(1)$ -approximation with high probability 64-approx

Our estimate \hat{J} satisfies $\frac{1}{C} \leq \hat{J} \leq Cd$ for some constant C .

with prob. $1 - \frac{\delta}{C}$ using $O(\log \log n)$ bits space for the max lsb and $O(\log n)$ bits for the hash function

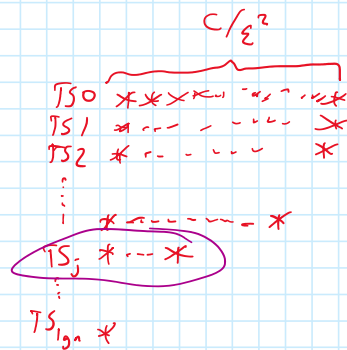
Refine to $(1 \pm \epsilon)$ solution

Trivial solution TS stores the first $\frac{C}{\epsilon^2}$ distinct elements.

TS is a $(1 \pm \epsilon)$ -solution if $d \leq \frac{C}{\epsilon^2}$.

Algorithm

1. Instantiate $TS_0, \dots, TS_{\lg n}$
2. Pick $g = [n] \rightarrow [n]$ from 2-wise family
3. Feed a_i to $TS_{\lfloor \lg(a_i) \rfloor}$
4. Output $B_j 2^{j+1}$ where $B_j = |TS_j| \approx \frac{1}{\epsilon^2}$



Let r.v. B_j be the # distinct elem hashed by g to TS_j .

$$\text{Let } \mathbb{E} B_j = \frac{d}{2^{j+1}} = Q_j.$$

By Chebyshev, $B_j = Q_j \pm O(\sqrt{Q_j})$ with good probability

$$= (1 \pm O(\epsilon)) Q_j \text{ if } Q_j \approx \frac{1}{\epsilon^2}$$

↑
 $O(\sqrt{Q_j}) \approx m(\pm)$

More precisely, use $j = \lfloor \lg \frac{\epsilon^2 d}{32} \rfloor$ where \uparrow is a constant approx.

$$= (1 \pm O(\epsilon)) / Q_j \quad \text{if } Q_j \approx \frac{1}{\epsilon^2}$$

$$\left[1 \pm \frac{O(\sqrt{Q_j})}{Q_j} \right] Q_j, \quad \frac{O(\frac{1}{\epsilon})}{\frac{1}{\epsilon^2}} = O(\epsilon)$$

Notice that Q_j 's decrease by factors of 2, and with good probability, one of the Q_j 's is close to $\frac{1}{\epsilon^2}$

Final space: $\frac{C}{\epsilon^2} \cdot (\lg n) \cdot (\lg n) \cdot (\lg \frac{1}{\delta}) = O(\epsilon^2 (\lg^2 n) (\lg \frac{1}{\delta}))$ bits

unique *strings entries + hash* *copies of Q_j* *medians to obtain high success prob.*

Can do better e.g. $O\left(\frac{1}{\epsilon^2} \lg \lg n + \underbrace{\lg n}_{\text{hash}}\right)$ for HyperLogLog

[Kane, Nelson, Woodruff, 2010, PODS]: $O\left(\frac{1}{\epsilon^2} + \lg n\right)$ optimal

F₂ second moment freq sketch

Consider a stream $S = \{a_1, \dots, a_n\}$, $a_i \in [m]$

Let $f_s = |\{a_i \mid a_i = s\}|$, the frequency of s in the stream.

We want to estimate $F_2 = \sum_{s=1}^m f_s^2$.

Algorithm:

- Let $h: [m] \rightarrow \{-1, 1\}$ be 4-wise independent.
- Let $x_s = h(s)$, Bernoulli r.v. with equal prob. of going to -1 or 1 .
- Output $a = \left(\sum_{i=1}^n h(a_i)\right)^2$ *i.e. each time we see s , add x_s to a variable, and square at the end.*

Notice: $a = \left(\sum_{i=1}^n x_s f_s\right)^2$

Claim: $\mathbb{E}(a^2) \leq 3(\mathbb{E}a)^2$

proof: $\mathbb{E}(a^2) = \mathbb{E}\left(\sum_{i=1}^n x_i f_i\right)^4 = \mathbb{E}\left(\sum x_s x_t x_u x_v f_s f_t f_u f_v\right)$

Urn $\perp (a) - \dots (a)$

proof. $E(a^2) = E\left(\sum_{s=1}^m x_s f_s\right)^2 = E\left(\sum_{1 \leq s, t, u, v \leq m} x_s x_t x_u x_v f_s f_t f_u f_v\right)$

Note, if any s, t, u, v are distinct, by 4-independence, the expectation is 0.

\Rightarrow only cases are either all 4 variables the same or 2 pairs

$$\begin{aligned} E(a^2) &= \binom{4}{2} E\left(\sum_{s=1}^m \sum_{t=s+1}^m x_s^2 x_t^2 f_s^2 f_t^2\right) + E\left(\sum_{s=1}^m x_s^4 f_s^4\right) \\ &= 6 \sum_{s=1}^m \sum_{t=s+1}^m f_s^2 f_t^2 + \sum_{s=1}^m f_s^4 \\ &\leq 3 \left(\sum_{s=1}^m f_s^2\right)^2 = 3 E^2 a. \end{aligned}$$

Thm 6.3 The average x of $r = \frac{2}{\epsilon^2 \delta}$ estimates a_1, \dots, a_r using independent sets of 4-way ind. hash functions h_s

$$\text{Prob}(|x - E x| > \epsilon E x) < \frac{\text{Var}(x)}{\epsilon^2 E^2 x} \leq \delta$$

proof. $\text{Var}(x) \leq \delta \epsilon^2 E^2 x$, and the rest follows by Chebyshev. \square