

Last class we saw a $(3, \delta)$ -approx. for the # of distinct elements in a data stream.

Setup: Data stream $S = \{s_1, \dots, s_m\}$

where m is HUGE

& for each i , $s_i \in \{1, 2, \dots, n\}$

Let $f = (f_1, \dots, f_n)$ where ~~$f_i = \# \text{ of } i \text{ in } S$~~

$$f_i = |\{j : 1 \leq j \leq m, s_j = i\}|$$

= # of occurrences of i in S

Our goal is to compute $Q = |\{i : f_i > 0\}|$

= # of distinct elements in S

Note, $Q = F_0$ where F_k was defined a couple of lectures ago.

Last class: We saw an approx. alg. that outputs \hat{Q} where:

$$\Pr\left(\frac{Q}{3} \leq \hat{Q} \leq 3Q\right) \geq .04$$

& by taking the median of $O(\log(1/\delta))$ trials we can boost this success prob. to $\geq 1 - \delta$.

this required space $O(\log(1/\delta) \log n)$.

②

Today: We'll boost the approximation factor:

for any $\epsilon > 0$, we'll guarantee that

$$\Pr(Q(1-\epsilon) \leq \hat{Q} \leq Q(1+\epsilon)) \geq 1-\delta.$$

& we'll do this with space $O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right) \log n\right)$.

Idea: Same basic approach as last time, which is to find max of $\text{zeros}(h(k))$ for $k \in S$, using a pairwise independent hash function.

Recall, for $k \in S$, $\Pr(\text{zeros}(h(k)) \geq r) = 2^{-r}$,

hence, we expect $\frac{Q}{2^Z}$ to have $\text{zeros}(h(k)) \geq Z$.

Our approach will be to keep track of how

many items have $\text{zeros}(h(k)) \geq Z$

where Z is close to $\max_k \text{zeros}(h(k))$.

Then, we'll output $|B|2^Z$ as our estimate of Q where B is the bucket containing those k with $\text{zeros}(h(k)) \geq Z$.

We'll keep $|B| \leq O(1/\epsilon^2)$ & if it exceeds that ~~it~~ then we'll increase Z .

Alg. [BJKST '04] = Bar-Yossef, Jayram, Kumar, Sivakumar & Trevisan. (3)

1. Choose a random pairwise independent hash function $h: [n] \rightarrow [n]$.
(We need n to be prime so if not choose $n < p \leq 2n$
which is prime)

Do this by choosing a, b indep. & uniformly from $\{0, 1, \dots, p-1\}$
& setting $h(i) = a + bi \pmod p$.

2. Set $z = 0$ & $B = \emptyset$.

3. Process data stream S in a one-by-one manner.

For element $k \in S$:

if $\text{zeros}(h(k)) \geq z$ then:

$B = B \cup (h(k), \text{zeros}(h(k)))$

while $|B| \geq c/\epsilon^2$: (specify c below)

$z = z + 1$

Remove all (α, β) from B
where $\beta < z$.

4. Output $|B|/2^z$

Recall from last class,

for $k \in \{1, \dots, n\}$, and integer $l \geq 0$,

$$\text{let } X_{l,k} = \begin{cases} 1 & \text{if } \text{zeros}(h(k)) \geq l \\ 0 & \text{o/w} \end{cases}$$

$$\& \text{ let } Y_l = \sum_{k: f_k > 0} X_{l,k} = |\{k \in S : \text{zeros}(h(k)) \geq l\}|$$

Consider the final value of z .

$$\text{The alg. outputs } \hat{Q} = Y_z 2^z$$

Let's bound the prob. that \hat{Q} is a poor approx. of Q :

we want to show that $\hat{Q} > (1+\epsilon)Q$ or $\hat{Q} < (1-\epsilon)Q$
are unlikely.

$$\begin{aligned} \hat{Q} > (1+\epsilon)Q &\Leftrightarrow Y_z 2^z > (1+\epsilon)Q \Leftrightarrow Y_z 2^z - Q > \epsilon Q \\ \& \hat{Q} < (1-\epsilon)Q &\Leftrightarrow Y_z 2^z < (1-\epsilon)Q \Leftrightarrow Y_z 2^z - Q < -\epsilon Q \end{aligned}$$

So, the alg. FAILS if $|Y_z 2^z - Q| \geq \epsilon Q$

$$\text{or } \left| Y_z - \frac{Q}{2^z} \right| \geq \frac{\epsilon Q}{2^z}$$

$$\text{Note, } E[Y_z] = \frac{Q}{2^z}$$

Note, $\Pr(X_{l,k}=1) = 2^{-l}$ & $E[X_{l,k}] = 2^{-l}$

Hence, $E[Y_l] = \frac{d}{2^l}$

and $\text{Var}(Y_l) = \sum_{k: f_k > 0} \text{Var}(X_{l,k}) \leq \sum_k E[X_{l,k}^2] = \sum_k E[X_{l,k}]$
 $\leq \frac{d}{2^l}$
 since $X_{l,k}$ is 0-1 r.v.

$\Pr(\text{FAIL}) = \Pr(|Y_z - \frac{d}{2^z}| > \frac{\epsilon d}{2^z})$

$= \sum_{r=1}^{\log n} \Pr(|Y_r - \frac{d}{2^r}| > \frac{\epsilon d}{2^r}, \& z=r)$

We'll choose s later $\leq \sum_{r=1}^{s-1} \Pr(|Y_r - \frac{d}{2^r}| > \frac{\epsilon d}{2^r}) + \sum_{r=s}^{\log n} \Pr(z=r)$

$= \sum_{r=1}^{s-1} \Pr(|Y_r - \frac{d}{2^r}| > \frac{\epsilon d}{2^r}) + \Pr(\# z \geq s)$

$= \sum_{r=1}^{s-1} \Pr(|Y_r - \frac{d}{2^r}| > \frac{\epsilon d}{2^r}) + \Pr(Y_{z-1} \geq \frac{c}{\epsilon^2})$

bound using Chebyshev's inequality

bound using Markov's inequality

because for z to \uparrow we need in the alg. that Y_{z-1} is too big.

(6)

$$\Pr\left(|Y_r - \frac{Q}{2^r}| \geq \frac{\epsilon Q}{2^r}\right) \leq \Pr\left(|Y_r - E[Y_r]| \geq \frac{\epsilon Q}{2^r}\right)$$

$$\leq \frac{\text{Var}(Y_r)}{\left(\frac{\epsilon Q}{2^r}\right)^2} \leq \frac{\left(\frac{Q}{2^r}\right)}{\left(\frac{\epsilon Q}{2^r}\right)^2} = \frac{2^r}{\epsilon^2 Q}$$

Note, $\sum_{r=1}^{s-1} 2^r \leq 2^s$

& hence $\sum_{r=1}^{s-1} \frac{2^r}{\epsilon^2 Q} \leq \frac{2^s}{\epsilon^2 Q}$

Now, choose s to be the largest integer where:

$$\frac{Q}{2^s} < \frac{24}{\epsilon^2}$$

note that by taking s one smaller $\frac{Q}{2^s}$ goes down by $\frac{1}{2}$

So we know $\frac{Q}{2^s} \geq \frac{12}{\epsilon^2}$

& note, $2^s \leq \frac{Q\epsilon^2}{12}$ i.e., $\frac{2^s}{Q} \leq \epsilon^2/12$.

So: $s = O(\log(cQ\epsilon^2))$ for some constant c .

⑦

$$\Pr(Y_{s-1} \geq \frac{c}{\epsilon^2}) \leq \frac{E[Y_{s-1}]}{c/\epsilon^2} \leq \frac{\epsilon^2 d}{c 2^{s-1}} \leq \frac{\epsilon^2 d}{c 2^s}$$

$$\leq \frac{2\epsilon^2 d}{c \epsilon^2}$$

since $\frac{d}{2^s} < \frac{24}{\epsilon^2}$

Therefore, combining these 2 bounds we have:

$$\Pr(\text{FAIL}) \leq \frac{2^s}{\epsilon^2 d} + \frac{2\epsilon^2 d}{c \epsilon^2}$$

$$\leq \frac{1}{\epsilon^2} \frac{\epsilon^2}{12} + \frac{48}{c} = \frac{1}{12} + \frac{48}{c} \leq \frac{1}{6}$$

for $c \geq 12.48$.

Now use the median of $O(\log(Y_s))$ indep't. trials to boost the success prob. to $\geq 1 - \delta$.

$$\begin{aligned}
 \text{Space: } & O(\log n) \times \frac{c}{\epsilon^2} + O(\log n) + O(\log \log n) \\
 & \begin{array}{cccc}
 \uparrow & \uparrow & \uparrow & \uparrow \\
 \text{per item} & \text{max \#} & \text{to choose } a, b & \text{for } z \leq \log n \\
 & \text{in } B & \text{store} &
 \end{array} \\
 & = O\left(\frac{1}{\epsilon^2} \log n\right)
 \end{aligned}$$

Can reduce it to $O\left(\log n + \frac{1}{\epsilon^2} (\log \frac{1}{\epsilon}) + \log \log n\right)$
 by using a hash function to
 keep track of the elements of B .