

Counting: Stream e_1, e_2, \dots, e_N

Goal: return an estimate of N using small space

• Morris Counter: $(1 \pm \epsilon, \frac{1}{2\epsilon^2})$ estimator

\hookrightarrow With prob. $1 - \frac{1}{2\epsilon^2}$: $|\hat{h} - N| \leq \epsilon \cdot N$

• Morris+ Counter: $\Theta(\frac{1}{\epsilon^2})$ independent instances
of Morris Counter

$\hookrightarrow (1 \pm \epsilon, \frac{1}{4})$

• Morris++ Counter: $\Theta(\log(\frac{1}{\delta}))$ independent instances of Morris+ Counter

$\hookrightarrow (1 \pm \epsilon, \delta)$

Median trick

$l = \lceil 12 \ln(\frac{1}{\delta}) \rceil$ independent instances of Morris+ Counter

return $\hat{h}_{\text{final}} = \text{median}(\hat{h}_1, \dots, \hat{h}_l)$

(\hat{h}_i : result of i -th instance)

Analysis: (We want: $(1-\epsilon) \cdot N \leq \hat{n}_{\text{final}} \leq (1+\epsilon) \cdot N$)
 otherwise: fails (to give correct approximation)

\hat{n}_{final} fails only if majority of $\hat{n}_1, \dots, \hat{n}_\ell$ fails

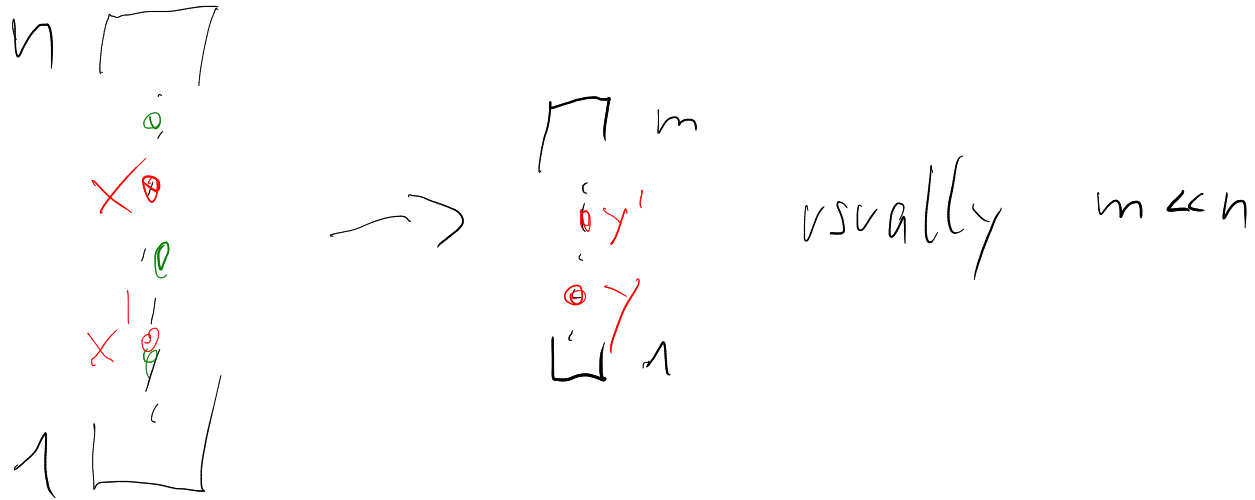
random variables $Z_i = \begin{cases} 1 & \text{if } \hat{n}_i \text{ fails} \\ 0 & \text{o.w.} \end{cases}$

By analysis of Morris counter: $\Pr[Z_i = 1] \leq \frac{1}{4}$ for all $1 \leq i \leq \ell$

Total number of failures $Z := \sum_{i=1}^{\ell} Z_i$ $\mu = p \cdot \ell = \frac{1}{4} \cdot \ell$

$$\Pr[\hat{n}_{\text{final}} \text{ fails}] = \Pr\left[Z \geq \frac{\ell}{2}\right] = \Pr\left[Z \geq 2 \cdot \frac{\ell}{4}\right] \leq \frac{1}{(1+\frac{1}{3})^\mu} \stackrel{\text{Chernoff Bound}}{\leq} \frac{1}{e^{\frac{1}{3} \cdot \mu}} = \frac{1}{e^{\frac{1}{12} \ell}} \leq \frac{1}{e^{\ln(\frac{1}{\delta})}} = \delta$$

Hash functions



Families of hash functions $h \sim \mathcal{H} \quad h: [n] \rightarrow m \quad [n] := \{1, 2, \dots, n\}$

• uniform $\forall x \in [n] \quad \forall y \in [m] \quad \Pr[h(x) = y] = \frac{1}{m}$

• universal $\forall x, x' \in [n] \quad \Pr[h(x) = h(x')] = \frac{1}{m}$
 \uparrow if $x \neq x'$

• 2-universal: $\forall x, x' \in [n] \quad \forall y, y' \in [m] \quad \text{if } x \neq x', \text{ then } \Pr[h(x) = y \wedge h(x') = y'] = \frac{1}{m^2}$

Remarks: • 2-universal \Rightarrow uniform

• 2-universal families of hash functions with small space and time consumption exist

Estimate # Distinct Element

... from a stream $e_1, \dots, e_N \in [h]$

frequency vector of stream $f = (f_1, \dots, f_n)$ where $f_j = |\{i : e_i = j\}|$

Goal: estimate $d := \{j \in [h] : f_j > 0\} =$

Think representing numbers in binary

zeros(p) = $\max \{i : 2^i \text{ divides } p\}$

= "number of 0s that p ends with"

$p = 110101000$

range $\overline{\hspace{10em}}$

end with 0 $\overline{\hspace{4em}}$

end with 00 $\overline{\hspace{2em}}$

end with 000 $\overline{\hspace{1em}}$

$O(\log n)$ space



Algorithm

Initialize: Choose random hash function $h : [h] \rightarrow [h]$ from 2-universal family
 $z \leftarrow 0$

Process token e_i : if $\text{zeros}(h(e_i)) > z$
 $z \leftarrow \text{zeros}(h(e_i))$

Output: $z^{2 + \frac{1}{2}}$

← at most $\log h$ zeros in any hash
 $O(\log \log h)$ space to store z
T.i.v., for final value of z

Analysis: for each $j \in [n]$, $r \geq 0$

random variables $X_{r,j} = \begin{cases} 1 & \text{if } \text{zeros}(h(j)) \geq r \\ 0 & \end{cases}$

$$Y_r = \sum_{j: f_j > 0} X_{r,j}$$

Thus: $Y_r > 0 \iff T \geq r$

$Y_r = 0 \iff T \leq r-1$

h uniform $j \rightarrow \frac{1}{n}$

fixed r, j

$$E[X_{r,j}] = \Pr[X_{r,j} = 1] = \Pr[\text{zeros}(h(j)) \geq r] = \frac{1}{2^r}$$

$$E[Y_r] = \sum_{j: f_j > 0} E[X_{r,j}] = \sum_{j: f_j > 0} \frac{1}{2^r} = \frac{d}{2^r}$$

By 2-universality $X_{r,j}$'s are independent

$$\text{Var}[Y_r] = \sum_{j: f_j > 0} \text{Var}[X_{r,j}] \leq \sum_{j: f_j > 0} E[X_{r,j}^2] = \sum_{j: f_j > 0} E[X_{r,j}] = \frac{d}{2^r}$$

$$\Pr[T \geq r] = \Pr[Y_r > 0] = \Pr[Y_r \geq \underbrace{1}_{\frac{1}{E[Y_r]}}] = \Pr\left[Y_r \geq \frac{2^r}{d} \cdot E[Y_r]\right] \leq \underbrace{\frac{d}{2^r}}_{\text{Markov Bound}}$$

$$\Pr[T < r] = \Pr[Y_r = 0] \leq \Pr\left[\underbrace{Y_r - \frac{d}{2^r}}_{E[Y_r]} \leq -\frac{d}{2^r}\right] \leq \frac{1}{\underbrace{\sqrt{\frac{d}{2^r}}}_{\sqrt{E[Y_r]}} \cdot \underbrace{\sqrt{\text{Var}[Y_r]}}_{\sqrt{\frac{d}{2^r}}}} = \frac{2^r}{d}$$

Let \hat{d} be estimate output by algorithm

$$\rightarrow \hat{d} = 2^{T+\frac{1}{2}} = \sqrt{2} \cdot 2^T$$

a : smallest integer s.t. $2^{a+\frac{1}{2}} \geq 3d$

$$\Pr[\hat{d} \geq 3d] = \Pr[T \geq a] \leq \frac{d}{2^a} = \frac{\sqrt{2}d}{2^{a+\frac{1}{2}}} \leq \frac{\sqrt{2}d}{3d} = \frac{\sqrt{2}}{3} < 0.5$$

b : largest integer s.t. $2^{b+\frac{1}{2}} \leq \frac{d}{3}$

$$\Pr[\hat{d} \leq \frac{d}{3}] = \Pr[T \leq b] = \Pr[T < b+1] \leq \frac{2^{b+1}}{d} = \frac{\sqrt{2} \cdot 2^{b+\frac{1}{2}}}{d} \leq \frac{\sqrt{2} \cdot \frac{d}{3}}{d} = \frac{\sqrt{2}}{3} < 0.5$$

$$\begin{aligned}
\Pr\left[\frac{d}{3} < \hat{d} < 3d\right] &= \Pr\left[\hat{d} > \frac{d}{3} \wedge \hat{d} < 3d\right] = 1 - \Pr\left[\hat{d} \leq \frac{d}{3} \vee \hat{d} \geq 3d\right] \\
&= 1 - \left(\Pr\left[\hat{d} \leq \frac{d}{3}\right] + \Pr\left[\hat{d} \geq 3d\right]\right) \\
&\leq 1 - 2 \cdot \frac{\sqrt{2}}{3} \approx 1 - \frac{2.828}{3} \geq 5\%
\end{aligned}$$

\hookrightarrow Boost this to success probability $1 - \delta$ by running $\Theta\left(\frac{1}{\delta}\right)$ instances of the algorithm and returning the median of the estimates (\rightarrow median trick)