

## Estimating Distinct Elements

$(1 \pm \epsilon)$  approx w. prob.  $(1 - \delta)$

Space  $O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$

$$h: U \rightarrow [0, 1] \quad Y = \min h(x) \quad E[Y] = \frac{1}{k+1}$$

$O\left(\log\left(\frac{1}{\delta}\right)\right)$  groups of  $O\left(\frac{1}{\epsilon^2}\right)$  hash fns each  
"median of means"

Sketching: "sketch" of data stream

Composable  $sk(S_1)$   $sk(S_2)$   $sk(S_1 \cup S_2)$

Assumption: hash fn  $u$  completely random

$$h(x) \in_R [0, 1]$$

$h(x_1)$   $h(x_2)$  ...  $h(x_k)$  independent

Pairwise Independent Hash fn

$H$ : family of hash fn  $h \in_R H$

$$P_h[h(x_1) = y_1, h(x_2) = y_2] = P_h[h(x_1) = y_1] \cdot P_h[h(x_2) = y_2]$$

$p$ : large prime  $x \in [p]$

$$h_{a,b}: [p] \rightarrow [p]$$

$$h_{a,b}(x) = ax + b \pmod{p}$$

$$H = \{h_{a,b}, a, b \in [p]\}$$

$$Y = \min \left\{ \frac{h(x)}{p} \right\}$$

$$P\left[Y < \frac{1}{3k}\right] < \frac{2}{5} \longrightarrow \in \left[0, \frac{1}{3k}\right] \quad k \cdot \frac{1}{3k} = \frac{1}{3}$$

$$P\left[Y > \frac{3}{k}\right] < \frac{1}{3} \longrightarrow \text{HW}$$

$\rightarrow O(\log(\frac{1}{\epsilon}))$  copies of hash fn & take median  
 constant then this is an estimate within  $[\frac{1}{3}, 3]$  of  $\frac{1}{K}$

$$a_i, b_i : h_i \quad a_i, b_i \in \mathbb{R}[p]$$

$(1+\epsilon)$  approx using pairwise independence  
 [Bar-Yossef et al, 2002]

Change algo: track smallest  $t$  hash elements

$$y_i : i^{\text{th}} \text{ smallest element} \quad E[y_i] = \frac{i}{K+1}$$

$$\text{Estimator: } \frac{t}{y_t} \approx K$$

Thm  $t = \frac{c}{\epsilon^2}$  with prob  $\geq \frac{2}{3}$   $\frac{(1-\epsilon)t}{K} \leq y_t \leq \frac{(1+\epsilon)t}{K}$

Pf (2nd ineq)

$$I = [0, (1+\epsilon)\frac{t}{K}]$$

$X_i$ : indicator for  $h(x_i) \in I$

$X = \sum X_i$  #hash values in  $I$

$$E[X] = \sum E[X_i] = K \cdot \frac{(1+\epsilon)t}{K} = (1+\epsilon)t$$

$$\begin{aligned} \Pr[y_t > \frac{(1+\epsilon)t}{K}] &= \Pr[X < t] \\ &= \Pr[X - E[X] < -\epsilon t] \\ &\leq \Pr[|X - E[X]| > \epsilon t] \end{aligned}$$

Chebyshev:  $\Pr[|X - E[X]| > \epsilon t] \leq \frac{\text{Var}(X)}{\epsilon^2 t^2}$

$$p = \Pr[X_i = 1] = \frac{(1+\epsilon)t}{K}$$

$$\begin{aligned}
 \mathbb{E}[X_i] &= p & \text{Var}(X_i) &= p(1-p) \\
 \mathbb{E}[X] &= kp & \text{Var}(X) &= kp(1-p) \leq \mathbb{E}[X] = (1+\epsilon)t \\
 & & & \text{uses pairwise independence}
 \end{aligned}$$

$$\Pr[|X - \mathbb{E}[X]| > \epsilon t] \leq \frac{(1+\epsilon)t}{\epsilon^2 t^2} = \frac{(1+\epsilon)}{\epsilon} \leq \frac{1}{6}$$

by suitable choice of  $c$



$O(\log(\frac{1}{\delta}))$  copies  
 Store  $t = O(\frac{1}{\epsilon^2})$  hash values  
 $(1+\epsilon)$  approximation w. prob.  $1-\delta$

Practical algo: HyperLogLog [Flajolet et al 2007]

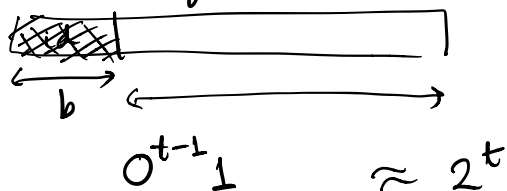
Assume: hash  $f^n$  is completely random  
 estimate cardinalities beyond  $10^9$

w. accuracy 2% using  $\sim 1.5$  Kbytes

Stochastic averaging: [Flajolet, Martin]

maintain  $m$  random variables  $m = 2^b$

break up stream into  $m$  substreams by using  
 first  $b$  bits of hash value



Each substream tracks max pos<sup>n</sup> of leading 1

$m(i)$ : pos<sup>n</sup> of leading 1

estimate  $2^{m(i)}$

Harmonic mean of these estimates

$$\sum_i 2^{-m(i)} \approx \frac{k}{m}$$

$$2^{-m(i)} \approx \frac{m}{k}$$

$$\sum_i 2^{-m(i)} \approx \frac{m^2}{k}$$

$$\text{Estimator} = \frac{m^2 \alpha_m}{\sum_i 2^{-m(i)}}$$

$$E[\text{Estimator}] \approx k$$

$$\frac{\sqrt{\text{Var}[\text{Estimator}]}}{k} \approx \frac{1.04}{\sqrt{m}}$$

Lower Bounds for Streaming Algorithms

Any deterministic algorithm that gives 1.4 approximation to # distinct elements must use  $\Omega(n)$  memory.

If we have inputs  $I_1, \dots, I_N$  for which algo must have distinct states, then  $\Omega(\log N)$  space.

$I_i, I_j$  & input  $I'$

$I_i \cup I'$  vs.  $I_j \cup I'$  have very different answers.

$\{S_i\}_{i=1}^N$  subsets of  $[n]$

$$\forall i \quad |S_i| = \frac{n}{10}$$

$$\forall i \neq j \quad |S_i \cap S_j| \leq \frac{n}{20}$$

$$S_i \cap S_j \quad |S_i \cup S_j| = \frac{n}{10}$$

$$|S_j \cup S_i| \geq \frac{3}{2} \cdot \frac{n}{10}$$

N different streams

$N = 2^{cn}$  many subsets. (Probabilistic method!)

$\log N = \Omega(n)$  lower bound.

Next Time: Frequency Moments.

$f_i$ : # of times that  $i$  appears

$$F_t = \sum_i f_i^t$$

$$F_2 = \sum_i f_i^2$$

$$F_1 = \sum_i f_i$$

[Alon, Matias, Szegedy 1996]

Streaming Model