

# Estimating Distinct Elements in a Data Stream

Among practical algorithms, 3 are state of the art

- Hyperloglog [Flajolet et al. 2002]
  - K-minimum values (KMV)
  - Adaptive Sampling
- space usage can all be made  
 $O(\log n + \frac{1}{\epsilon^2} \cdot (\log(\frac{1}{\epsilon}) + \log \log n))$   
for  $\delta = \frac{1}{50}(\epsilon)$ .

I will describe ~~them~~ <sup>KMV + Adaptive Sampling</sup> as having space usage  $O(\frac{\log n}{\epsilon^2})$ . I will explain why shortly.

In theory: [KMV] gave an  $O(\log n + \frac{1}{\epsilon})$  space algorithm with  $O(1)$  update time.

~~Let~~ KMV: Let  $h: [n] \rightarrow [0,1]$  be a random function.

- Track the  $K$ -smallest hash values observed in the stream. Let  $m_k$  denote  $k$ 'th smallest hash value seen.
- Output the estimate  $\frac{k-1}{m_k}$ .

Adaptive sampling! Let  $h: [n] \rightarrow [0,1]$  be a random hash function.

- Initialize  $i \leftarrow 0$ .
- While processing stream, store all hash values  $< 2^{-i}$ .
- If more than  $K$  hash values are stored, set  $i \leftarrow i+1$ .
- If  $S$  is set of hash values stored at end of stream, output  $\frac{|S|}{2^i}$ .

Facts about KMV and Adaptive Sampling.

- KMV has a natural min-heap based implementation. However, this requires  $O(\log K)$  time per stream update, and doubles its space usage compared to just keeping the  $K$  hash values in a hash table.
- ~~Adaptive Sampling is faster (0(1)-amortized update time) and doesn't oscillate due to periodic purges.~~  
Adaptive Sampling is faster (0(1)-amortized update time) and doesn't oscillate due to periodic purges.

- Both KMV and Adaptive Sampling are unbiased, <sup>we will prove this for KMV</sup>
- $\text{Var}[KMV] \leq \frac{F_0^2}{K-1}$
- $\text{Var}[Adaptive Sampling] \approx \frac{1.44 F_0^2}{K-1}$

We will not prove either of these facts, ~~here~~ we note that ~~the means~~ by Chebyshev, the variance bounds mean setting  $K = O(\frac{1}{\epsilon^2})$  is enough for an  $(\epsilon F_0)$ -multiplicative approximation.

i.e.  $\Pr[|KMV - F_0| > \epsilon F_0] \leq \frac{1}{4}$ .

$O(KMV)$  with  $K = \frac{1}{\epsilon^2} + 1$  counters  $\approx$  about  $\frac{F_0}{\sqrt{4\epsilon^2}} = \frac{\epsilon}{2} \cdot F_0$

• What is the space usage of KMV and adaptive sampling?

~~It is enough for the hash function families that form a pairwise independent hash family [Ch]~~

• It turns out it is enough for the hash function  $h$  to be from a pairwise independent hash family mapping  $[n]$  to a set of size  $O(\frac{1}{\epsilon^2} \log^2 n)$ . Hence,  $h$  requires  $O(\log n + \log(\frac{1}{\epsilon}))$  bits to represent, and each of the  $O(\frac{1}{\epsilon^2})$  hash values take  $O(\log(\frac{1}{\epsilon}) + \log \log n)$  bits to represent, for a total space bound of  $O(\log n + \frac{1}{\epsilon^2} \cdot (\log(\frac{1}{\epsilon}) + \log \log n))$ .

• In practice, it can be very useful to store not just the hash values but also the corresponding identifiers. If you do this, the space usage is  $O(\frac{\log n}{\epsilon^2})$ . This lets you, e.g., obtain estimates for ad hoc subsets of users, where the subset of interest is only determined at query time.

Hyperloglog:

- Let  $h$  be a random function mapping  $[n]$  to  $[0, 1]$
- Let  $g$  be a random function mapping  $[n]$  to  $[k]$
- For each  $j \in [k]$ , track  $v_j := \max_{i: g(a_i)=j} \text{zeros}(h(a_i))$
- Let  $Z$  be the harmonic mean of  $2^{v_j}$  values

Output  $\propto k \cdot Z$ , ( $k$  is a constant meant to correct a small bias in  $k \cdot Z$ )

Intuition:  $2^{v_j}$  should be about  $\frac{F_0}{k}$ , so  $k \cdot Z$  should be about  $F_0$ .

To get an  $(\epsilon, \delta)$ -multiplicative approximation, can take  $k \approx \frac{1}{\epsilon^2}$ , so total space usage is  $O(\log n + \frac{\log \log n}{\epsilon^2})$ .

Proof that KMV is unbiased. For any stream  $\sigma$ ,

Fix any  $j \in [n]$  appearing one or more times in the stream.

Define  $V_j := \begin{cases} \frac{1}{m_k} & \text{if } h(j) \text{ is one of the } k \text{ smallest hash values observed} \\ 0 & \text{otherwise} \end{cases}$

Let  $H_{-j}$  denote the  $F_0 - 1$  hash values of all other items in the stream other than  $j$ .

Let  $m_k$  denote the  $k$ 'th smallest hash value in  $H_{-j}$ , and  $m_{k-1}$  the  $(k-1)$ 'st smallest.

$$\text{Then } \mathbb{E}[V_j | H_{-j}] = \Pr_h[h(j) < m_{k-1}] \cdot \frac{1}{m_{k-1}}$$

$$+ \Pr[m_{k-1} < h(j)] \cdot 0$$

$$= \frac{m_{k-1}}{m_{k-1}} = 1$$

Note the KMV estimate is  $\sum_j V_j$ , so  $\mathbb{E}[\text{KMV estimate}] = \sum_{j \in \text{stream}} \mathbb{E}[V_j] = \sum_{j: \text{seen}} 1 = F_0$

Analysis of Adaptive Sampling with  $K = \frac{c}{\epsilon^2}$  ( $c$  to be specified later)

Let  $X_{r,j} = \begin{cases} 1 & \text{if } h(j) \leq 2^{-i_r} \\ 0 & \text{otherwise} \end{cases}$

and let  $Y_r = \sum_{j: f_j > 0} X_{r,j}$ . Let  $i_{\max}$  denote the value of  $i$  at the end  
set of hash values less than  $2^{-i_{\max}}$

So the returned estimate  $\hat{F}_0 = \frac{1}{2^{i_{\max}}} \cdot Y_{i_{\max}}$

Exactly as in last lecture, we obtain

Fact 1:  $\mathbb{E}[Y_r] = \frac{F_0}{2^r}$ ,  $\text{Var}[Y_r] \leq \frac{F_0}{2^r}$ . Fact 1!

Note: if  $i_{\max} = 0$  then the algorithm never saw more than  $K$  hash values  
 so it returns exactly  $|S| = F_0$ .

otherwise, we need to bound the probability  $|Y_{i_{\max}} \cdot 2^{i_{\max}} - F_0| \geq \epsilon F_0$   
 $\Leftrightarrow |Y_{i_{\max}} - \frac{F_0}{2^{i_{\max}}}| \geq \frac{\epsilon F_0}{2^{i_{\max}}}$ . call this event FAIL.

Let  $s$  be the unique integer s.t.  $\frac{12}{\epsilon^2} \leq \frac{F_0}{2^s} \leq \frac{24}{\epsilon^2}$ .

Then  $\Pr[\text{FAIL}] = \sum_{r=1}^{\infty} \Pr\left[|Y_r - \frac{F_0}{2^r}| \geq \frac{\epsilon F_0}{2^r} \text{ and } i_{\max} = r\right]$   
 $\leq \left(\sum_{r=1}^{\infty} \Pr\left[|Y_r - \frac{F_0}{2^r}| \geq \frac{\epsilon F_0}{2^r}\right]\right) + \left(\sum_{r \geq s_0} \Pr[i_{\max} = r]\right)$   
 $\leq \dots + \Pr[i_{\max} \geq s_0]$   
 $+ \Pr[Y_{s-1} \geq \frac{K}{\epsilon^2}]$

By Chebyshev and Fact 2's Markov's

$$\Pr\left[\left|Y_r - \frac{F_0}{2^r}\right| \geq \frac{\epsilon F_0}{2^r}\right] \leq \frac{2^r}{\epsilon^2 F_0}$$

So the first term is at most  $\frac{\epsilon^2}{12}$

$$\sum_{r=1}^{s-1} \frac{2^r}{\epsilon^2 F_0} \leq \frac{2^s}{\epsilon^2 F_0} \leq \frac{1}{12}$$

By Markov's inequality and Fact 1, the second term is at most

$$\frac{\mathbb{E}[Y_{s-1}]}{\frac{\epsilon}{2}} = \frac{F_0}{2^{s-1}} \cdot \frac{\epsilon^2}{c} \leq \frac{48}{c} \leq \frac{1}{12} \text{ if we choose } c=4$$

$\uparrow \leq \frac{48}{\epsilon^2}$

In total,  $\Pr[\text{Fail}] \leq \frac{1}{12} + \frac{1}{12} = \frac{1}{6}$  so there is an  $(6, \frac{1}{6})$ -approximation algorithm