

6.854/18.419: Advanced Algorithms

Last Time: consistent hashing and random trees

↑
hash functions that
evolve well

↑
routing schemes that
work with inconsistent news

Today: Distinct elements and count-min

Meta Question: what can we do if we can't store
our data, but it is streaming by?

$O(\log n)$
bits of
memory

Problem #1: Count the number of distinct elements
in a sequence

x_1, x_2, x_3, \dots

How many distinct words did Shakespeare use?

31,534 of which 14,376 only used once

"Using only memory equivalent to 5 lines of
printed text, you can estimate with a typical
accuracy of 5% and in a single pass
[the number of words Shakespeare used]"

Durand and Flajolet 2003

Attempt #1: Choose a random hash function

$$h: U \rightarrow [0, 1]$$

Pass through data $h(x_1), h(x_2), h(x_3), \dots$ and compute min

Let $Y = \min(h(x_1), h(x_2), h(x_3), \dots)$

Lemma: If there are N distinct elements in

then $\mathbb{E}[Y] = \frac{1}{N+1}$

Pf:

$$\begin{aligned} \mathbb{E}[Y] &= \int_0^1 \Pr[Y \geq z] dz && \text{(convince yourself in discrete case)} \\ &= \int_0^1 N z (1-z)^{N-1} dz \\ \text{cdf} \Rightarrow \text{pdf} &= \int_0^1 (1-z)^N dz = \left. -\frac{(1-z)^{N+1}}{N+1} \right|_{z=0}^{z=1} \\ &= \frac{1}{N+1} \quad \square \end{aligned}$$

Alternatively $\mathbb{E}[Y] =$ "probability of choosing $N+1$ values in $[0,1]$, and last is min"

$$= \frac{1}{N+1} \quad (\text{by symmetry})$$

How do we bound Y 's deviation? Markov? Chebyshev? Chernoff?

Lemma: $\text{var}(Y) < \left(\frac{1}{N+1}\right)^2$

Pf:

$$\mathbb{E}[Y^2] = \int_0^1 N z^2 (1-z)^{N-1} dz = \frac{2}{(N+1)(N+2)}$$

$$\text{var}(Y) = \underbrace{E[Y^2]} - (\underbrace{E[Y]})^2$$

$$\frac{2}{(N+1)(N+2)} - \left(\frac{1}{N+1}\right)^2 = \left(\frac{1}{N+1}\right) \left(\frac{N}{(N+1)(N+2)}\right) \quad \square$$

Alternatively it is a beta distribution with $\alpha=1, \beta=N$

$$\text{var} = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$$

But the standard deviation is comparable to expectation
(why is this bad?)

any ideas how to fix?

Attempt #2 choose k hash functions almost Flajolet-Martin

$$h_1, h_2, \dots, h_k: U \rightarrow [0, 1]$$

Find minimum Y_1, Y_2, \dots, Y_k for each one

$$\text{Set } \bar{Y} = \frac{1}{k} \sum_{i=1}^k Y_i, \text{ output } 1/\bar{Y} - 1$$

Analysis: $E[\bar{Y}] = 1/(N+1), \text{var}(\bar{Y}) < \frac{1}{k(N+1)^2}$

Thus by Chebyshev

$$\Pr\left[|\bar{Y} - \frac{1}{N+1}| > \frac{\epsilon}{N+1}\right] \leq \frac{\frac{1}{k(N+1)^2}}{\frac{\epsilon^2}{(N+1)^2}} = \frac{1}{k\epsilon^2}$$

Then $\frac{N+1}{1+\epsilon} \leq \frac{1}{\bar{Y}} \leq \frac{N+1}{1-\epsilon}$, what about bit complexity of min?

(First, Misra-Gries)

Problem #2: ϵ -approximate Frequency Counts

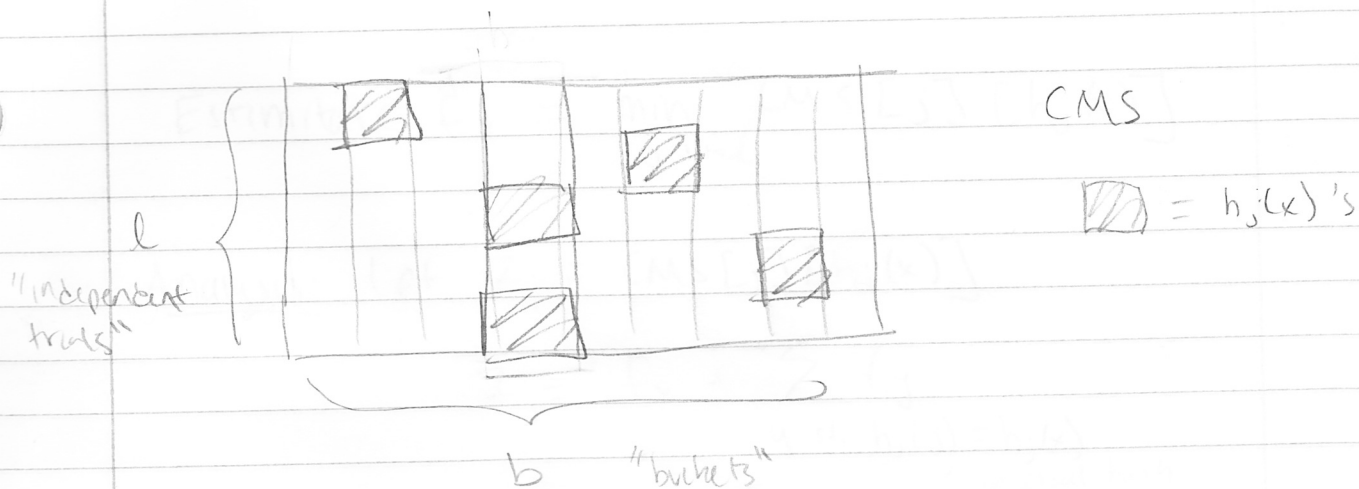
Given a sequence $x_1, x_2, x_3, \dots, x_n$ output a list of values s so that

(1) every value that occurs at least $\frac{n}{k}$ times is on list

(2) no value that occurs less than $\frac{n}{k} - \epsilon n$ times is on list

be able to answer queries
 $\text{count}(x) \leq \epsilon n$
weaker goal:
answer queries on frequencies

Count-Min Sketch: Cormode, Muthukrishnan 2005



Choose l random hash functions $h_1, h_2, \dots, h_l: (U \rightarrow [b])$

For all x ,

For $j = 1$ to l

Increment $\text{CMS}[j][h_j(x)]$

How do we approximate the number of times x appeared?

Let $f_x =$ frequency of x

claim: For all x, j

$$\text{CMS}[j][h_j(x)] \geq f_x$$

PF: Fix j . Each time we see x , we increment

$$\text{CMS}[j][h_j(x)]$$

although we can increment it for other reasons (collisions) \square

$$\text{Estimate: } \hat{f}_x = \min_{j=1 \dots d} \text{CMS}[j][h_j(x)]$$

Analysis: Let $Z_j = \text{CMS}[j][h_j(x)]$

$$Z_j = f_x + \sum_{y \text{ st. } h_j(y) = h_j(x)} f_y$$

$$\text{Hence } \mathbb{E}[Z_j] = f_x + \frac{\sum_y f_y}{b} \leq f_x + \frac{n}{b}$$

$\sum_y f_y = n$ (universal hash)

$$\text{Set } b = \frac{2}{\epsilon}, \text{ then } \Pr[Z_j - f_x \geq \epsilon n] \leq \frac{1}{2}$$

markov $Z_j - f_x$ nonnegative

$$\text{Finally } \Pr\left[\min_{j=1 \dots d} Z_j \geq f_x + \epsilon n\right] \leq \frac{1}{2^d}$$

\uparrow
 $O(\log n)$

Space: $O(k \log n)$ where $\epsilon = \frac{1}{k}$

heavy hitters

Frequent Items, Misra-Gries 1982

Given a sequence $x_1, x_2, x_3, \dots, x_n$ output a list with

(1) at most k value

(2) every value that appears at least $\frac{n}{k+1} + 1$ times is on list

Initialize empty list

For each item

If same as some item on list, increment its counter

Else if less than k items in list, add to list set its counter to one

Else decrement all counters, delete items whose counter reaches zero

let f_x = frequency of value x

Lemma: At end, value x 's counter is at least $f_x - \frac{n}{k+1}$

either x is decremented or not added

Pf: Each elimination of x eliminates k other symbols

Hence no more than $\frac{n}{k+1}$ occurrences of x can be eliminated. \square

space: $O(k \log n)$