

Notes/Reminders

- webpage: <https://onak.pl/ds563>
cs543
- HW 0 / HW 1 (deadline extended to 2/1)
- electronic device policy

Last time:

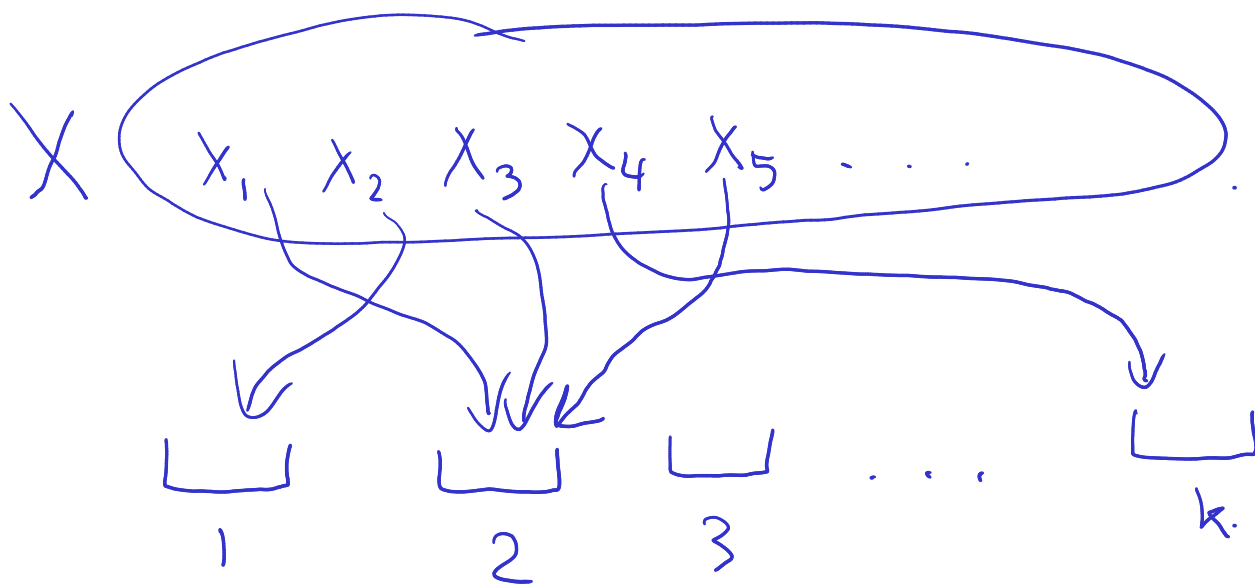
- Problem: provide estimates of what fraction of items in an evolving collection is any given element
- Started: CountMin Sketch

Today:

- Wrap up CountMin Sketch (CMS)
 - Important concepts: Linear Sketches and Streaming Algorithms
 - ~~- Deterministic estimates (Misra-Gries)~~
- Next time

CMS Last Time

Randomly map elements of universe X
to buckets via $h: X \rightarrow [k]$
" $\{1, 2, \dots, k\}$



- keep count $A[i]$ of number of elements in each bucket i
- fraction estimate for $y \in X$:

$$\frac{A[h(y)]}{\sum_{i=1}^k A[i]} \stackrel{\text{def}}{=} S \leftarrow \text{total number of items}$$

For a fixed y :

$$\frac{f(y)}{s} \leq \text{estimate} \leq \frac{f(y)}{s} + \frac{2}{k}$$

\uparrow always true
 \uparrow w.p. $\geq 1/2$
 Set $k = \lceil 2/\epsilon \rceil$ to overestimate by at most ϵ w.p. $\geq 1/2$

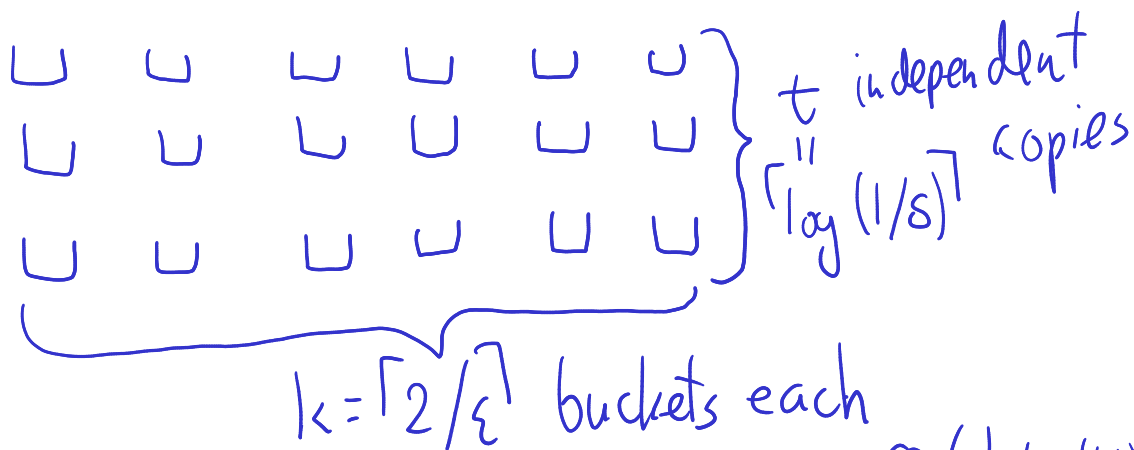
$f(x) = \#$ occurrences of x

How to make probability of overestimating by more than ϵ at most $s \in (0, 1/2)$?

- Run $t = \lceil \log(1/s) \rceil$ independent copies
- Query y : return the minimum of estimates from all copies

$$\Pr \left[\text{all overestimate by more than } \epsilon \right] \leq \left(\frac{1}{2} \right)^t \leq s$$

Visually:



3-3

Total space: $O\left(\frac{1}{\epsilon} \log(1/s)\right)$

What is missing?

How do we store fully random hash functions?

- We can't! Lots of space!

- Pairwise independence suffices

For $x \neq y$: $\mathbb{E}[C_{x,y}] = \Pr[h(x)=h(y)] \leq \frac{1}{k}$

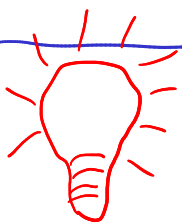
or even

$\mathbb{E}[C_{x,y}] \leq \frac{O(1)}{k} \leftarrow$ fixed constant

(just increase k by a constant factor)

- discussion section:

universal hashing / simple construction for strings



Linear sketches

(CountMin sketch is an example)

maintained sketch \rightarrow $\begin{bmatrix} \\ \\ \end{bmatrix}$

For CMS, all the k buckets

$\begin{bmatrix} 3 & -4 \end{bmatrix}$ $q \times 1$

$\begin{bmatrix} \text{randomized} \\ \text{matrix} \end{bmatrix}$

what our algorithm does $q \times N$

$\begin{bmatrix} \\ \\ \end{bmatrix}$

input: frequency vector $f(\dots)$

$N \times 1$ dimensions

Linear sketching algorithm projects high-dimensional input onto low dimensional sketch

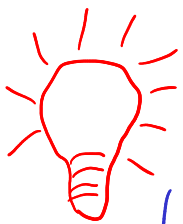
Nice properties of linear sketches:

- Can handle both insertions & deletions (insertion & deletion are the "same" change to the sketch but with the opposite sign)
- Can be computed separately for subsets and easily combined via addition

$$\left[\quad \right] (f_1 + f_2 + \dots + f_n) = \left[\quad \right] f_1 + \dots + \left[\quad \right] f_n$$

Application: different centers

handling different parts of the data set



Streaming Algorithms

Streaming algorithm A

huge stream of data

← 7, 3, 8, 11, 563, 47, ...

- A reads and processes items one by one
- A should use much less space than the input size

A few comments:

- Items in the stream can be anything

Common examples: integers, points from
some metric space, database records,
graph edges, ...

- CMS is a streaming algorithm

- Here the algorithm gets to read
the stream once ("one pass"),

but sometimes a small number
of passes over the stream makes

sense in some scenarios, e.g.,

data living on a storage device

in which case sequential reading

can maximize throughput