

Today

- Wrap up CountMin Sketch
- Heavy hitters
- Second moment estimation via AMS sketch

CountMin Sketch Recap

Goal:

- items arrive in arbitrary order (from some set X)
- multiple copies ~~of~~ allowed
- provide estimates "what fraction of items is $y \in X$?"
- use small space

Last time:

- $h: X \rightarrow \{1, \dots, k\} = [k]$ random hash function

- keep array $A[1..k]$

- initially $A[i] = 0$ for all $i \in [k]$

- insertion of element $x \in X$:

$$A[h(x)] \leftarrow A[h(x)] + 1$$

- estimate for $y \in X$:

$$g(y) = \frac{A[h(y)]}{\sum_{i=1}^k A[i]}$$

Properties:

$f(x)$ = exact number of occurrences of x

$$S = \sum_{i=1}^k A[i] = \sum_{x \in X} f(x)$$

Always: $\frac{f(y)}{S} \leq g(y)$

With probability $1/2$: $g(y) \leq \frac{f(y)}{S} + \frac{2}{k}$

To get additive ϵ approximation, set $k = \lceil 2/\epsilon \rceil$

How to make probability of error at most $\delta \in (0, 1/2)$?

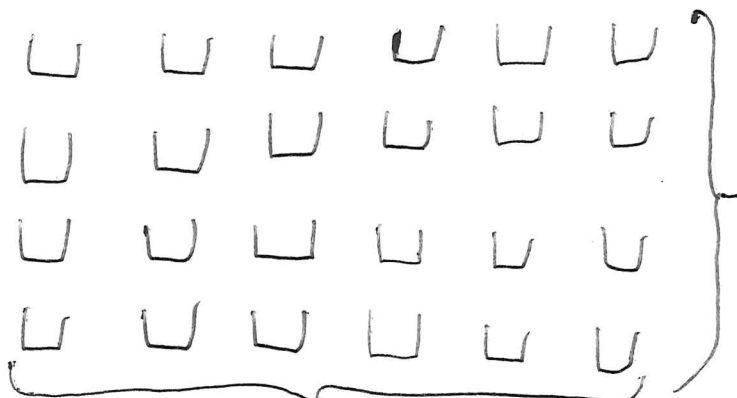
- Run $t = \lceil \log(1/\delta) \rceil$ independent copies

- Query y : return the minimum of estimates from all copies

$$\Pr[\text{all wrong}] \leq \left(\frac{1}{2}\right)^t \leq \delta$$

||
overestimate by more than ϵ

Visually:



t rows

||
 $\lceil \log(1/\delta) \rceil$

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

$k = \lceil 2/\epsilon \rceil$ buckets each

What is missing?

How do we store random hash functions?

- We can't (lots of space!)
- Pairwise independence suffices

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

- Or even

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

(just increase k by a constant factor)

- discussion section: sample construction for strings

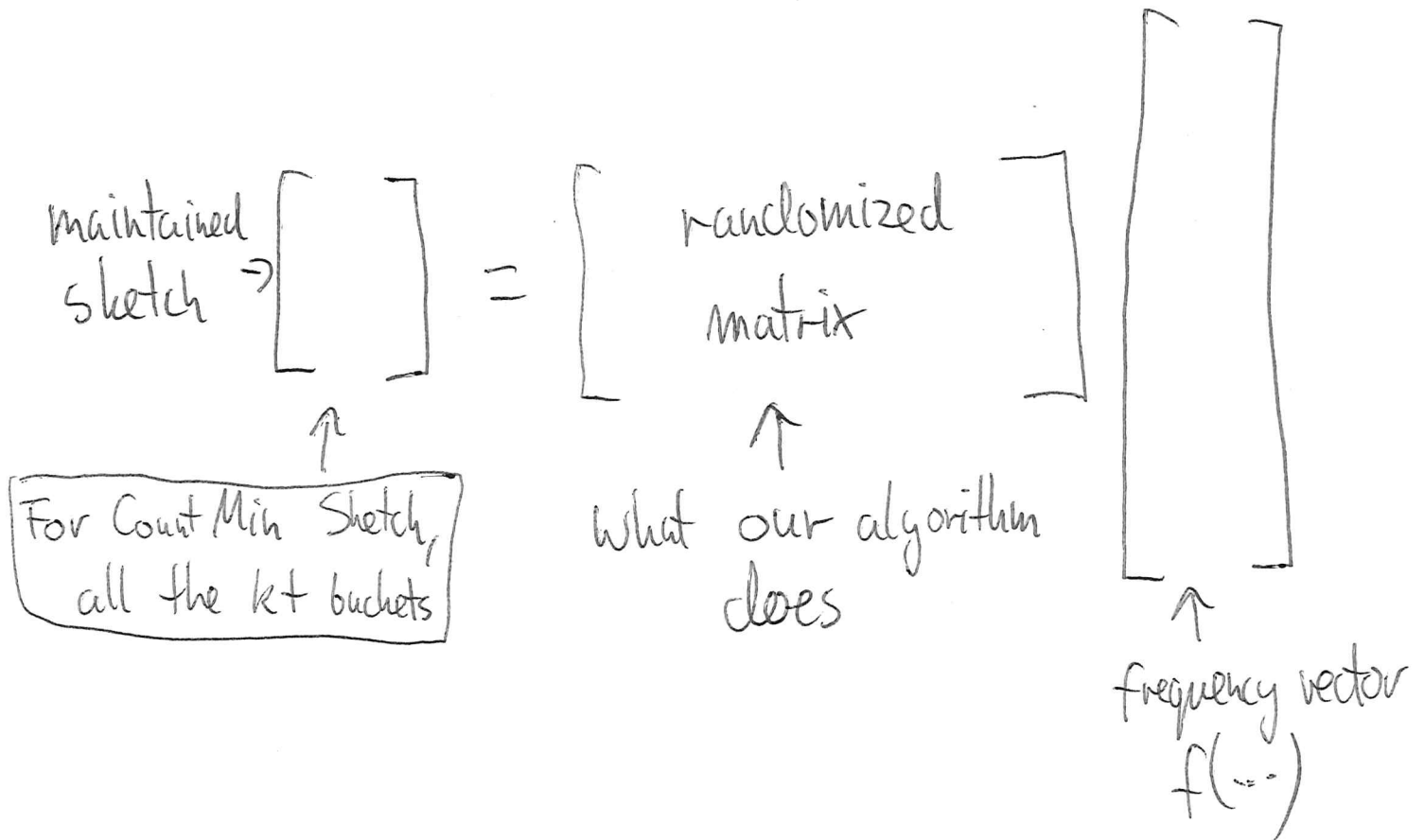
Nice properties of CountMin Sketch:

- can handle deletions (decrease corresponding counters)
- can be computed separately for subsets and easily combined (use the same hash functions!)

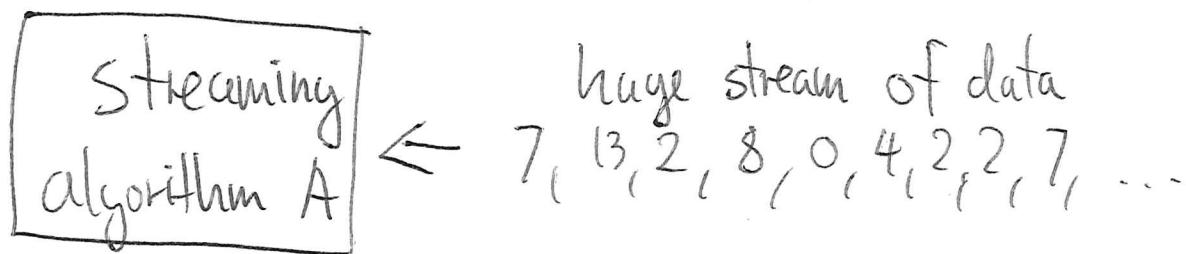
Application: different data centers handling different parts of the data set

Important concept: linear sketches

Count Min Sketch is an example



Important concept: streaming algorithms



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

Heavy hitters: list frequent elements in the stream

Setting: same as for CountMin Sketch

X - universe from which elements of the stream come from

$f(x)$ - number of occurrences of $x \in X$

S - total number of items

Our task: For some $\epsilon \in (0, 1)$, ^{parameter}

return $H \subseteq X$ s.t.

$$\forall x \in X: \begin{aligned} f(x) \geq 2\epsilon \cdot S &\Rightarrow x \in H \\ f(x) \leq \epsilon S &\Rightarrow x \notin H \end{aligned}$$

Approach:

- find candidates $H' \subseteq X$

- use Count Min Sketch to verify:

output all $x \in H'$ for which

CM Sketch says "fraction of $x \geq 2\epsilon$ "

How to find small H' ?

TO BE CONTINUED