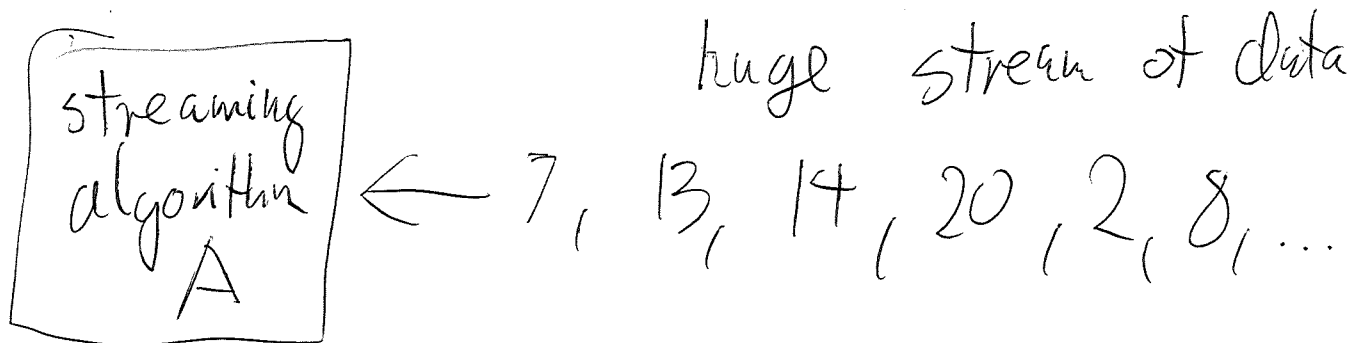


DS 563, Fall 2021, Lecture 2: Heavy Hitters

Streaming algorithms



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

Today: Heavy hitters (i.e. list frequent elements in the stream)

~~Setup~~: Setting: same as for Count Min Sketch

X - universe of items that we get to see

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

S - total number of items

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

return $H \subseteq X$ s.t.

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

Approach: find candidates, use Count Min sketch to verify

$$\uparrow \\ H' \subseteq X$$

output all $x \in H'$ for which Count Min sketch says ~~they~~ fraction of occurrences $\geq 2\epsilon$

How to find H' ?

Warm-up: find element that occurs (called leader)
more than $1/2$ of the time
or output nothing or any $z \in X$
if there is no leader

Algorithm:

- remember at most one element plus a count (= number of copies)

- when new item x arrives

- if storage empty ~~at that~~ store x with a count of 1

- otherwise, if same item in storage, increase ~~the~~ the count
if different items, throw away the new item + one copy of the item in storage (= decrease

→ the count)
If the count becomes 0, empty storage

Why works:

if x is a leader in \forall set $Y \subseteq X$,
it will still be ~~a leader~~
a leader after removing two
different elements

If there is a leader,
it will be in the storage
at the end of the stream

Easy extension to finding elements
more frequent than $1/k$ fraction

- store at most k elements with counts
- if k different elements in storage,
remove a copy of each

Homework: additional properties + estimate
space usage