# Lecture 2: Sketching

Streaming algorithms: Setting where data appears sequentially, and the goal is to process this in an online fashion, using few resources
$\hookrightarrow$ usually space.

Problem 1: Given a stream of numbers, track majority if it exists.

$n$ numbers in stream

$$X_1, \ldots, X_n \longrightarrow \log n \text{ bit } \#s.$$

$$m = maj(X_1, \ldots, X_n) \text{ if } \#\{i : X_i = m\} > \frac{n}{2}.$$

naive: $n \log n$ bits of memory.

Claim: Can do with $O(\log n)$!

Algo:

count = 0, guess = NULL.

For $X_i$ in stream:

    If count = 0,

        current = $X_i$, count = 1.

    else if $X_i$ = guess

        count ++

    else

        count --

output guess

2  1  4  2  4  4  1  4  2  2  2  2  2  2

Question: how to analyze?
    Hint: consider signed counter

Heavy hitters problem: Given $X_1, \ldots, X_n$
say that $y$ is an HH if
$$\#\{i : X_i = y\} \geq n/k \quad \leftarrow \text{ some parameter}$$
$$k = n/2 + 1$$

Hard for even reasonable $k$!

Problem: Output the most common element
of the stream (the "heaviest hitter").

Claim: Requires $\Omega(n \log n)$ space!

$$00 \quad X_1, X_2 \cdots X_n \quad \Big| \quad \begin{matrix} x \in S? \\ x \notin S \end{matrix} \qquad X_i \in \{1, \ldots, n^2\}$$

$$\longrightarrow \binom{n^2}{n} \quad \text{such sequences.}$$

If I had $\leq B$ bits of
memory, $\leq 2^B$ distinct
states. But each distinct
state has different answer.

useful approx:
when $k \ll n$,
$$\binom{n}{k} \approx n^k$$

$$2^B \geq \binom{n^2}{n} \approx n^{2n}$$

$$B \geq n \log n.$$

For $k \leq n/2$, do same except repeat last
element many times!

( Breaks only at majority ).

$$\Omega(n \log n) \text{ for all "reasonable" } k.$$

---

Relax : $\varepsilon - HH.$ params $k, \varepsilon$.

Given stream $x_1, \cdots, x_n$:

1). If $x$ occurs $\geq n/k$ times, output it

2). If we output $x$, then it occurs $\geq n/k - \varepsilon n$ times.
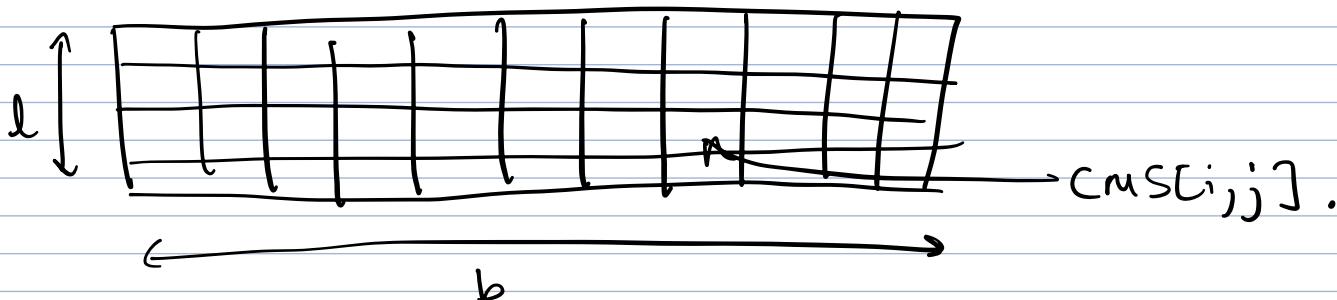
Space: $O(1/\varepsilon)$

e.g. if $\varepsilon = 1/2k$: output all $x \geq n/k$

any output $x \geq \frac{n}{2k}$.

using $O(k)$ memory

## Tool: Count-Min Sketch

Implements 2 operations: $Inc(x)$, $Count(x)$.

Params $\quad b = \#$ buckets $\qquad (b \approx 1/\varepsilon)$

$\qquad \quad \ell = \#$ hash functions. $\quad (\ell \approx O(1))$



$\to CMS[i,j]$.

$h_1, \cdots, h_\ell : \mathcal{U} \to \{0, \cdots, b-1\}$.

are your "nice" hash functions.

$Inc(x)$: $CMS[i, h_i(x)] {+}{+} \quad \forall i = 1, \cdots, \ell$.

$Count(x)$: $\min\limits_{i} CMS[i, h_i(x)]$

why does this work? Let $x$ occur $C_x$ times in stream.

Know: $Count(x) \geq C_x$ (why?).

need to bound $\underline{overestimation}$

let

$$Z_i = CMS[i, h_i(x)] = C_x + \underbrace{\sum\limits_{\substack{y \neq x \\ h_i(y) = h_i(x)}} C_y}_{} = (*)$$

$\forall x \neq y, \quad Pr[h_i(y) = h_i(x)] = \frac{1}{b}.$

$\Rightarrow \quad \mathbb{E}[C_*] = \sum_{y \neq x} \mathbb{E}\left[\mathbb{1}[h_i(y) = h_i(x)]\right] = \frac{n - C_x}{b}$

$\leq \frac{n}{b}$

Set $\quad b = \frac{2}{\varepsilon}.$

$\Rightarrow \quad \leq \frac{\varepsilon n}{2}$

Markov's Inequality: $\quad Y > 0$ is r.v.,

$$P[Y \geq \alpha] \leq \frac{\mathbb{E}[Y]}{\alpha}.$$

$\Rightarrow Pr[C_* \geq \varepsilon n] \leq \frac{\varepsilon n / 2}{\varepsilon n} \leq \frac{1}{2}.$

$\Rightarrow \quad Pr[Z_i - C_x \geq \varepsilon n] \leq \frac{1}{2}.$

$Pr[\min(Z_1, \cdots, Z_\ell) \geq C_x + \varepsilon n] \leq \left(\frac{1}{2}\right)^\ell$

$\ell = \log_2 \frac{1}{\delta}. \longrightarrow \leq \delta.$

$\Rightarrow$ For $\text{connt}(x)$ to be accurate to $\varepsilon n$ w.p. $1 - \delta$,

need to set $\quad b = \frac{2}{\varepsilon}, \quad \ell = \log \frac{1}{\delta}$

space $= O\left(\frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$ words of memory.