

Streaming

- Model
- Reservoir Sampling
- Counting Distinct Elements
- Hash families (pairwise independence)
- Heavy Hitters

(β) - Heavy Hitters

INPUT: A stream $s_1 \dots s_n \in \{1..N\}$

GOAL: List of all elements that appear
 $\geq \beta \cdot n$ times

$$\beta = 1/2$$

DATA STRUCTURE: "COUNT-MIN-SKETCH"

- List all heavy hitters
- Given any $a \in \{1..N\}$, let $f_a = \text{frequency of 'a'}$
 $f_a \pm \epsilon n$

$M[i, j] = \# \text{ of elements in the stream that map to Bucket } j \text{ under } h_i$

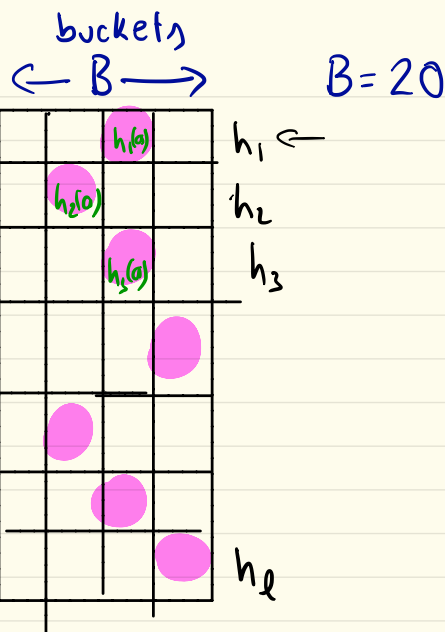
INITIALIZE:

- $M[i, j] = 0$

- Hash functions

$$h_1 \dots h_\ell : \{1 \dots N\} \rightarrow \underbrace{\{1 \dots B\}}_{\text{buckets}}$$

$$\ell \log n = l$$



INCREMENT (element a)

for $i = 1$ to ℓ

$$M[i, h_i(a)]++$$

QUERY (element a) \approx approximately estimate
frequency of a .

$$\text{RETURN } \min \left\{ \begin{array}{l} M[1, h_1(a)] \\ M[2, h_2(a)] \\ \vdots \\ M[l, h_l(a)] \end{array} \right\} \geq f_a$$

Fix i , and $a \in \{1..N\}$

F

$$\mathbb{E}_{h_i} M[i, h_i(a)] = \mathbb{E}_{h_i} \left(f_a + \sum_{\substack{h_i(b)=h_i(a) \\ b}} f_b \right)$$

pairwise independence

$$= f_a + \left(\sum_b f_b \right) \left(\frac{1}{B} \right)$$

For a fixed $b \neq a$

$\Pr_{h_i} [h_i(b) = h_i(a)] = 1/B$

$$\leq f_a + \frac{1}{B} (n) = f_a + \frac{n}{B}$$

MARKOV'S INEQUALITY: $X \geq 0$ be a r.v.

$$\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$$

$$\Pr[M[i, h_i(a)] - f_a \geq 2n/B] \leq 1/2$$

$$\Pr \left[\min_i \{M(i, h_i(\alpha))\} > f_\alpha + 2n/B \right] \quad \left[\begin{array}{l} B = 2/\epsilon \\ \Rightarrow 2n/B = \lceil \epsilon n \rceil \end{array} \right]$$

$$= \Pr \left[(M(1, h_1(\alpha)) > f_\alpha + 2n/B) \text{ And } (M(2, h_2(\alpha)) \geq f_\alpha + 2n/B) \right. \\ \left. \text{And } (M(l, h_l(\alpha)) \geq f_\alpha + 2n/B) \right]$$

(h_1, h_2, \dots, h_l are chosen independently)

$$= \Pr [M(1, h_1(\alpha)) > f_\alpha + 2n/B] \cdot \Pr [\quad] \cdots \Pr [M(l, h_l(\alpha)) \geq f_\alpha + 2n/B]$$

$$\leq \frac{1}{2} \cdot \frac{1}{2} \cdots \frac{1}{2} = \frac{1}{2^l} = \frac{1}{2^{10 \log n}} \approx \frac{1}{n^{10}}$$

β -Heavy Hitters:

s_1, \dots, s_n

for $i = 1$ to n

increment (s_i)

if ($\text{query}(s_i) \geq \beta n$) Add s_i to
list of β -heavy hitters.