

Count-Min Sketch

Counting (frequent) elements from stream $e_1, \dots, e_N \in [n]$ $e_i = (i_t, \Delta_t)$
add $\Delta_t \in \mathbb{Z}$ to f_{i_t}

Algorithm Table C k rows w columns
specified later

Let h_1, \dots, h_k be 2 -universal independent hash functions from $[n] \rightarrow [w]$

For each item $e_i = (i_t, \Delta_t)$ in the stream do

For $l=1$ to k do

$$C[l, h_l(i_t)] \leftarrow C[l, h_l(i_t)] + \Delta_t$$

For each $i \in [n]$ set

$$\tilde{f}_i = \min_{l \in [k]} C[l, h_l(i)]$$

15? 10

	1	2	3	4	5	6
1	10	20	1	2	5	8
4	4	7	9	1	7	1
8	8	1	7	2	3	4

Theorem If always $f \geq 0$ and $k = \Omega(\log \frac{1}{\delta})$ and $w > \frac{2}{\epsilon}$, then for any fixed $i \in [n]$: $\Pr[\tilde{f}_i \geq f_i + \epsilon \|f\|_1] \leq \delta$

$$\|f\|_1 = \sum_{i \in [n]} |f_i| \quad \text{1-norm}$$

Proof: Fix $i \in [L]$

Random variable $Z_\ell = \mathbb{C}[h_\ell(i)]$: value of counter in row ℓ that i is hashed to
 $\forall \ell \in [k]$

$$\begin{aligned} \mathbb{E}[Z_\ell] &= f_i + \sum_{i \neq j} f_j \cdot \Pr[h_\ell(j) = h_\ell(i)] = f_i + \sum_{j \neq i} \frac{1}{w} f_j \\ &= f_i + \frac{1}{w} \sum_{j \neq i} f_j \leq f_i + \frac{\varepsilon}{2} \cdot \|f\|_1 \end{aligned}$$

$\underbrace{\sum_{j \neq i} f_j}_{\leq \|f\|_1}$ \uparrow 2-universal h_j

$$\mathbb{E}[Z_\ell - f_i] \leq \frac{\varepsilon}{2} \cdot \|f\|_1$$

By Markov Inequality, $\Pr[Z_\ell - f_i \geq \varepsilon \|f\|_1] \leq \frac{1}{2}$

$$\Pr\left[\min_{\ell \in [k]} Z_\ell \geq f_i + \varepsilon \|f\|_1\right] = \Pr\left[\bigwedge_{\ell \in [k]} Z_\ell \geq f_i + \varepsilon \|f\|_1\right]$$

$$\leq \left(\frac{1}{2}\right)^k \leq \delta$$

□

Typical: set $\delta = \frac{1}{n^{c+1}}$ (for some constant $c \geq 1$)

$$\rightarrow \Pr[\tilde{f}_i > f_i + \varepsilon \|f\|_1] \leq \frac{1}{n^{c+1}} \text{ for each } i$$

\Rightarrow For all i simultaneously

$$\Pr[\forall i: \tilde{f}_i \leq f_i + \varepsilon \|f\|_1]$$

$$= 1 - \Pr[\exists i: \tilde{f}_i > f_i + \varepsilon \|f\|_1]$$

$$\leq n \cdot \frac{1}{n^{c+1}} = \frac{1}{n^c} \text{ Union Bound}$$

$$= 1 - \frac{1}{n^c} \text{ ("with high probability")}$$

$$\text{Space: } O(\underbrace{k \cdot w \cdot \log \|f\|_1}_{\text{counters}} + \underbrace{k \cdot \log n}_{\text{hash functions}}) =$$

$$= O\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} \cdot (\log n + \log \|f\|_1)\right)$$

Sketch

n-dimensional

sketch is a summary of \forall data D s.t.

$$\text{sketch}(D_1 \circ D_2) = \text{sketch}(D_1) + \text{sketch}(D_2)$$

composability property

In particular, linear sketches are of the form $\text{sketch}(D) = \Pi \cdot D$

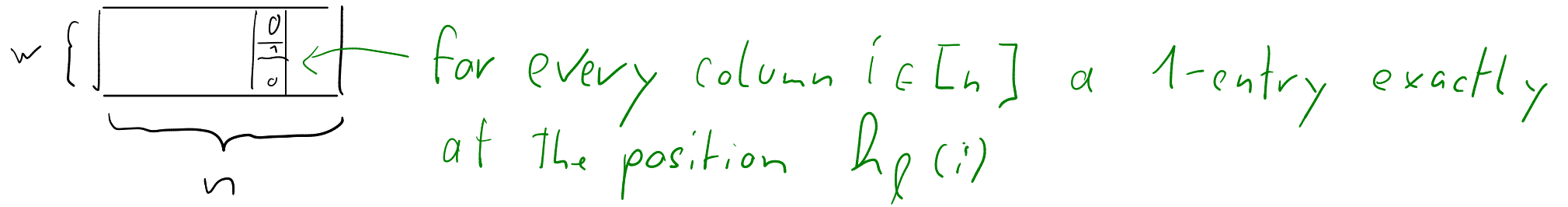
\uparrow
 $t \times n$ matrix

CountMin as a Linear Sketch

Need to show that sketch vector produced by CountMin sketch (for fixed hash functions) can be obtained by multiplying a fixed sketch matrix by the stream's frequency vector. (But we do not do this explicitly in the algorithm, only in the analysis)

Specifically, for a fixed hash function $h_\ell: [n] \rightarrow [w]$

we can define $w \times n$ matrix Π_ℓ



$\rightarrow \Pi_l f \rightarrow w$ -dimensional vector representing the counters of the hashed items

$$\Pi = \begin{pmatrix} \pi_1 \\ \pi_2 \\ \vdots \\ \pi_k \end{pmatrix}$$