

Algorithms for (strict) turnstile streams  $\sigma = \langle (a_1, d_1), \dots, (a_m, d_m) \rangle$

Point queries: Goal is to output a sketch from which one can derive, for any  $i \in [n]$ , an estimate  $\hat{f}_i$  of  $f_i$  such that

$$(*) \quad 0 \leq \hat{f}_i - f_i \leq \epsilon \cdot \|f\|_1$$

$\uparrow$   
 $= \sum_i f_i = M.$

Subtle distinction in randomized case:

- "For-all error guarantee": ~~with~~ with probability  $\geq 1 - \delta$ , (\*) hold simultaneously for all estimates  $\hat{f}_i$  returned by the algorithm.
- "For-each error guarantee": For each  $i \in [n]$ , (\*) holds with probability  $1 - \delta$ .

One can turn any For-each sketch into a for-all sketch with an  $O(\log n)$  blowup in space by reducing failure probability to  $\frac{\delta}{n}$  using the median trick and then union bounding over all  $i$ .

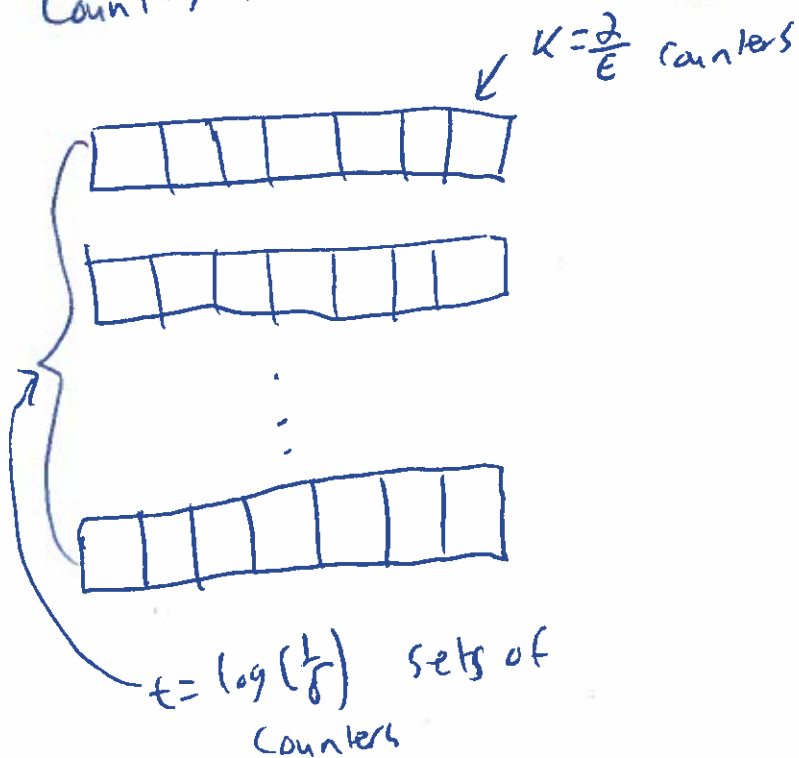
Recall: • *Mira-Gries* achieves this goal for insert-only streams using  $O(\frac{\log n}{\epsilon})$  bits of space (since it is deterministic, it is automatically a for-all sketch).

• There is a trivial way to turn *Mira-Gries* into a turnstile streaming algorithm, but the error will grow like

$$\epsilon \cdot \sum_{j=1}^m |\delta_j| \text{ instead of } \epsilon \cdot \sum_{i=1}^n f_i.$$

Run separate instances of *Mira-Gries* on positive increments and negative increment updates. Let  $f_i$  be the difference of the estimates returned for  $i$  by the two instances. Error is at most the sum of the errors in the two estimates so at most  $\epsilon \cdot \sum_j |\delta_j|$ .

### Count-Min Sketch [CM05]



• Choose  $t$  hash functions  $h_1, \dots, h_t : [n] \rightarrow [k]$  at random from a pairwise independent hash family

• When processing update  $(a_j, \delta_j)$ :

• For  $l = 1 \dots t$

$$C[l][h_l(a_j)] += \delta_j$$

• On query output  $\hat{f}_i := \min_{1 \leq l \leq t} C[l][h_l(i)]$

In the strict turnstile model, it is clear each  $\hat{f}_i$  is always an overestimate of  $f_i$ , since a counter's value is just the sum of the frequencies of all items that hash to it.

Claim: For any fixed  $i \in [n]$ ,  $0 \leq \hat{f}_i - f_i \leq \epsilon \cdot M$  with probability  $\geq 1 - \delta$ .  
This is a per-counter error guarantee.

Proof: We analyze the "excess" in each counter  $C[l][h_\ell(i)]$  that  $i$  hashes to. Let  $X_\ell := C[l][h_\ell(i)] - f_i$  denote this excess. For each  $j \neq i$ , let

$$Y_{\ell,j} = \begin{cases} 1 & \text{if } h_\ell(i) = h_\ell(j) \\ 0 & \text{otherwise} \end{cases} \quad \text{Then } X_\ell = \sum_{j \in [n] \setminus \{i\}} f_j \cdot Y_{\ell,j}.$$

By pairwise independence of the hash family,  $\mathbb{E}[Y_{\ell,j}] = \frac{1}{K}$ . Thus, by linearity of expectation,

$$\mathbb{E}[X_\ell] = \sum_{j \in [n] \setminus \{i\}} \frac{f_j}{K} = \frac{M - f_i}{K}$$

Since each  $f_j \geq 0$ ,  $X_\ell$  is a non-negative random variable,  $\int_0^\infty \mathbb{1}_{x \geq t} dx = t$  we can apply Markov's inequality to conclude  $\Pr[X_\ell \geq \epsilon \cdot M] \leq \frac{1}{\epsilon \cdot K} \stackrel{\epsilon \cdot K \geq 2}{\leq} \frac{1}{2}$ .

Since the hash functions are mutually independent,

$$\begin{aligned} \Pr[\hat{f}_i - f_i \geq \epsilon \cdot M] &= \Pr[X_\ell \geq \epsilon \cdot M \text{ for all } \ell \in [t]] \\ &= \prod_{\ell=1}^t \Pr[X_\ell \geq \epsilon \cdot M] \leq \left(\frac{1}{2}\right)^t \leq \delta. \quad \blacksquare \end{aligned}$$

Note: To get a formal error guarantee, need to increase space usage from

$$O\left(\frac{1}{\epsilon} \log\left(\frac{1}{\delta}\right) (\log m + \log n)\right) \text{ to } O\left(\frac{1}{\epsilon} \cdot \log\left(\frac{n}{\delta}\right) (\log m + \log n)\right).$$

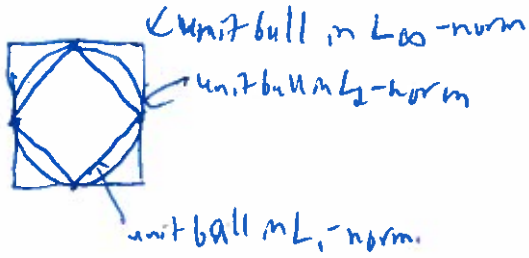
# Count sketch [CLP04]

Guarantee of this sketch: For each fixed  $i \in [n]$ , one can derive an estimate  $\hat{f}_i$  of  $f_i$  such that with probability  $1-\delta$ ,  $|f_i - \hat{f}_i| \leq \epsilon \cdot \|f\|_2$ .

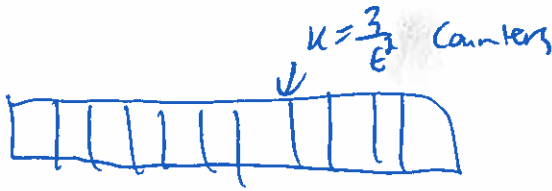
The sketch uses space  $O(\frac{1}{\epsilon^2} \cdot \log(\frac{1}{\delta}) \cdot (\log m + \log n))$ .

Note: For any vector  $f$ ,  $\|f\|_1 \leq \|f\|_2$ , so it is incomparable to Counting.

It uses a factor  $\frac{1}{\epsilon}$  more space, but its error might be smaller.



Algorithm:



Full sketch runs  $\log(\frac{1}{\delta})$  copies of basic estimator and outputs the median estimate.

**Basic Estimator**

- Choose a <sup>random</sup> hash function  $h: [n] \rightarrow [k]$  from a pairwise independent hash family.
- Choose a random hash function  $g: [n] \rightarrow \{-1, 1\}$  from a pairwise independent hash family.
- When processing update  $(a_j, f_j)$ :  

$$C[h(a_j)] \pm f_j \cdot g(a_j)$$
- on query  $i$ , output  $\hat{f}_i := g(i) \cdot C[h(i)]$

Analysis of Basic Estimator: Fix  $i \in [n]$  and let  $X := \hat{f}_i$ . <sup>Claim 1:  $\mathbb{E}[X] = f_i$</sup>  For each  $j \neq i$ , let

$$Y_j := \begin{cases} f_j & \text{if } h(j) = h(i) \\ 0 & \text{otherwise} \end{cases} \quad \text{Then } X = g(i) \cdot \sum_{j=1}^n f_j \cdot g(j) \cdot Y_j = f_i + \sum_{j \in [n] \setminus \{i\}} f_j \cdot g(i) \cdot g(j) \cdot Y_j$$

Since  $g, h$  are independent we have:

$$\mathbb{E}[g(i) \cdot g(j) \cdot Y_j] = \mathbb{E}[g(i)] \cdot \mathbb{E}[g(j)] \cdot \mathbb{E}[Y_j] \quad \text{by pairwise indep. of } g$$

$$\stackrel{(*)}{=} 0 \cdot \mathbb{E}[Y_j] = 0$$

So by linearity of expectation,

$$\mathbb{E}[X] = f_i + \sum_{j \in [n] \setminus \{i\}} f_j \cdot \mathbb{E}[g(i) \cdot g(j) \cdot Y_j] \stackrel{by (*)}{=} f_i + \sum_{j \in [n] \setminus \{i\}} f_j \cdot 0 = f_i$$

Similar to (\*) we also have:

$$\begin{aligned} (**) \text{ For any } j \neq j', \quad \mathbb{E}[g(i) \cdot g(i') \cdot Y_j \cdot Y_{j'}] &= \mathbb{E}[g(i)] \cdot \mathbb{E}[g(i')] \\ &\quad \cdot \mathbb{E}[Y_j \cdot Y_{j'}] \\ &= 0 \cdot 0 \cdot \mathbb{E}[Y_j \cdot Y_{j'}] = 0 \end{aligned}$$

In addition we have:

$$(***) \quad \mathbb{E}[Y_j^2] = \mathbb{E}[Y_j] = \Pr[h(i) = h(i)] = \frac{1}{k}.$$

Claim 2:  $\text{Var}[X] \leq \frac{\|f\|_2^2}{K}$ .

Proof:  $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \mathbb{E}[X^2] - f_i^2$ .

$$\mathbb{E}[X^2] = \mathbb{E}\left[\left(f_i + \sum_{j \in [n] \setminus \{i\}} f_j \cdot g(i, j) \cdot Y_j\right)^2\right]$$

distributive law

$$\downarrow = \mathbb{E}\left[f_i^2 + 2f_i \sum_{j \in [n] \setminus \{i\}} f_j \cdot g(i, j) \cdot Y_j + \left(\sum_{j \in [n] \setminus \{i\}} f_j \cdot g(i, j) \cdot Y_j\right)^2\right]$$

linearity of expectation

$$\downarrow = f_i^2 + 2f_i \underbrace{\sum_{j \in [n] \setminus \{i\}} f_j \mathbb{E}[g(i, j) \cdot Y_j]}_{= 0 \text{ by } (*)} + \mathbb{E}\left[\left(\sum_{j \in [n] \setminus \{i\}} f_j \cdot g(i, j) \cdot Y_j\right)^2\right]$$

linearity of expectation

$$\downarrow = f_i^2 + \mathbb{E}\left[\sum_{j \in [n] \setminus \{i\}} f_j^2 \cdot \underbrace{g(i, j)^2}_{=1} \cdot \underbrace{Y_j^2}_{=1} + \sum_{j \neq j' \in [n] \setminus \{i\}} f_j \cdot f_{j'} \cdot \underbrace{g(i, j) \cdot g(i, j')}_{=1} \cdot Y_j \cdot Y_{j'}\right]$$

$$= f_i^2 + \mathbb{E}\left(\sum_{j \in [n] \setminus \{i\}} f_j^2 Y_j^2 + \sum_{j \neq j' \in [n] \setminus \{i\}} f_j \cdot f_{j'} \cdot g(j, j') \cdot Y_j \cdot Y_{j'}\right)$$

linearity of expectation

$$\downarrow = f_i^2 + \sum_{j \in [n] \setminus \{i\}} f_j^2 \mathbb{E}[Y_j^2] + \sum_{j \neq j' \in [n] \setminus \{i\}} f_j \cdot f_{j'} \cdot \mathbb{E}[g(j, j') \cdot Y_j \cdot Y_{j'}]$$

by (\*\*\*)

$$\downarrow = f_i^2 + \sum_{j \in [n] \setminus \{i\}} f_j^2 K$$

$= 0 \text{ by } (***)$

Hence  $\text{Var}[X] = \mathbb{E}[X^2] - f_i^2 = \cancel{f_i^2} + \sum_{j \neq i} f_j^2 / k - \cancel{f_i^2}$

$$\leq \frac{\|F\|_2^2}{k} \quad \blacksquare$$

By Chebyshev,  $\Pr[|X - \mathbb{E}[X]| \geq \underbrace{\epsilon \cdot \|F\|_2}_{= \epsilon k \cdot \text{Var}[X]}] \leq \frac{1}{(\epsilon \cdot k)^2} = \frac{1}{\epsilon^2 k^2} \stackrel{\text{choose } k}{\leq} \frac{1}{3}$ .

Since the basic estimator gives an estimate with error  $\leq \epsilon \cdot \|F\|_2$  with probability  $\geq 2/3$ , the median of  $O(\log(1/\delta))$  basic estimates satisfies the error bound w.p.  $\geq 1 - \delta$ .