

# COMMONWEALTH OF AUSTRALIA

## Copyright Regulations 1969

### WARNING

This material has been reproduced and communicated to you by or on behalf of the University of Sydney pursuant to Part VB of the Copyright Act 1968 (**the Act**). The material in this communication may be subject to copyright under the Act. Any further copying or communication of this material by you may be the subject of copyright protection under the Act.

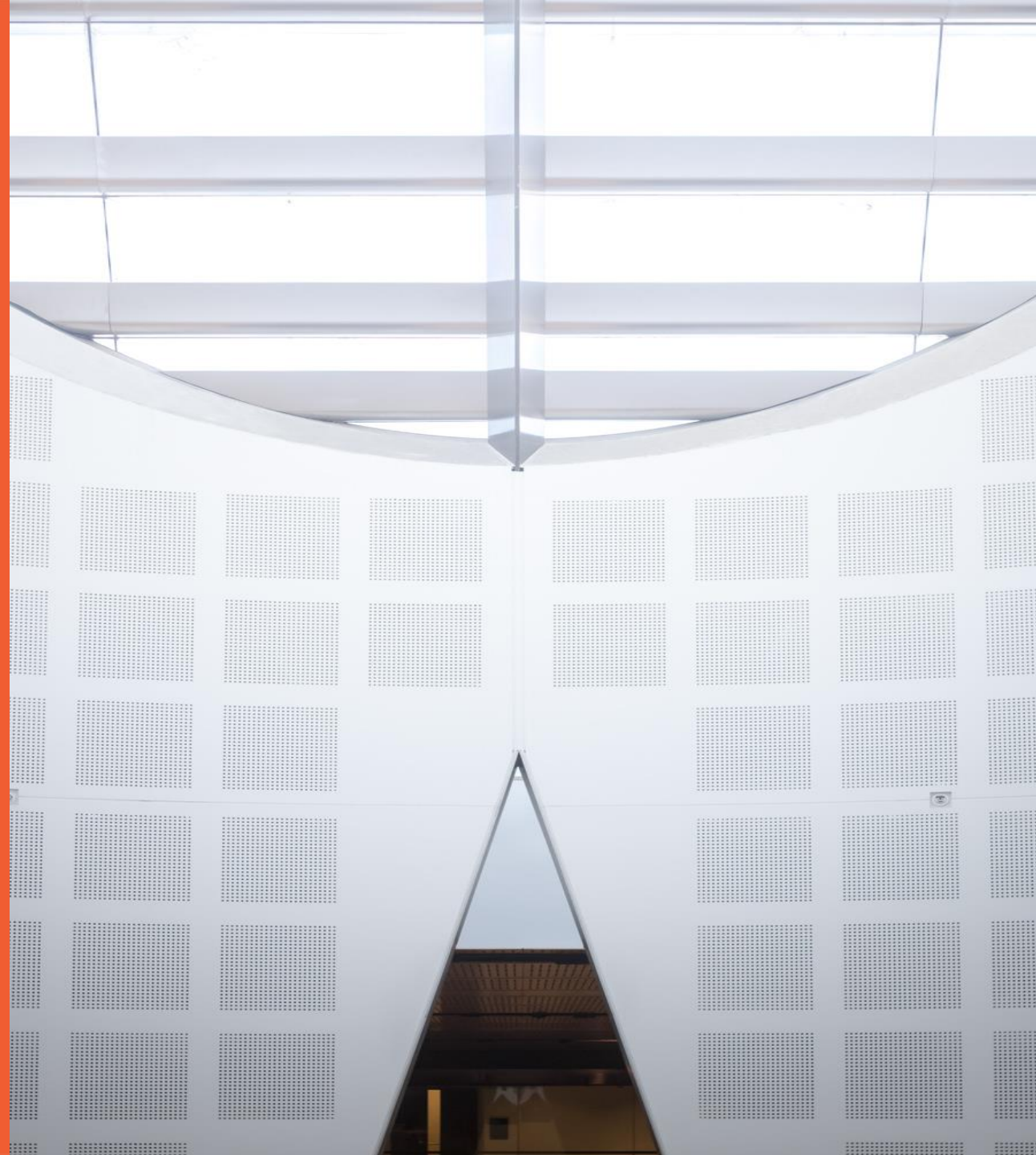
**Do not remove this notice.**

COMPx270: Randomised and  
Advanced Algorithms  
Lecture 9: Streaming and  
Sketching II

Clément Canonne  
School of Computer Science



THE UNIVERSITY OF  
SYDNEY



## A question

You design a streaming algorithm **A** to solve some problem. The stream of data arrives.

## A question

You design a streaming algorithm **A** to solve some problem. The stream of data arrives.

628, 516, 163, 509, 15, 499, 772, 588, 737, 439, 79, 866, 186, 18,  
854, 459, 146, 518, 748, 737, 685, 188, 939, 724, 27, 719, 263, 795,  
120, 573, 853, 132, 522, 3, 298, 123, 932, 993, 180, 674, 1, 619, 989,  
142, 496, 178, 191, 524, 716, 501, 677, 712, 452, 768, 591, 551, 439,  
397, 229, 214, 43, 639, 353, 610, 737, 203, 933, 279, 877, 30, 513,  
518, 616, 714, 633, 804, 422, 731, 867, 184, 124, 881, 595, 193, 254,  
240, 4, 260, 303, 319, 757, 723, 309, 365, 278, 512, 658, 233, 393,  
875

## A question

You design a streaming algorithm **A** to solve some problem. The stream of data arrives. **A** outputs its answer:

21.5

## A question

Oh, no! That wasn't the end. More data arrives.

## A question

Oh, no! That wasn't the end. More data arrives.

834, 992, 528, 12, 181, 274, 159, 150, 716, 71, 755, 4, 324, 398, 802,  
176, 302, 941, 678, 934, 546, 753, 812, 47, 755, 721, 893, 53, 410

## A question

Oh, no! That wasn't the end. More data arrives. **A** outputs its answer on this:

18.1



## A question

How do you combine 21.5 and 18.1 to get the answer on the whole data stream?

## Sketching

$\sigma \xrightarrow{A} S(\sigma)$  : sketch  
(succinct "summary"  
of the stream  $\sigma$ )  
Problem-dependent

$\sigma' \xrightarrow{A} S(\sigma')$

[From  $S(\sigma)$ ,  
can  
get  
answer  
to problem]

$f(S(\sigma), S(\sigma')) \approx S(\sigma \circ \sigma')$   
some way to combine things  
(concatenation of  $\sigma, \sigma'$ )

## Even better: linear sketching



## Even better: linear sketching

" $b = +$ "

$$S(s_1) + S(s_2) = S(s_1 + s_2)$$

↑ what is addition?

Typically, addition of vectors in  $\mathbb{R}^k$   
(and sketches are vectors:  $S(s) \in \mathbb{R}^k$ )

Equivalent:

$$S: \mathbb{R}^n \rightarrow \mathbb{R}^k \quad k \ll n$$

$$\begin{array}{ccc} k \downarrow & \left( \begin{array}{c} S \end{array} \right) & \\ & \xleftarrow{n} & \left( \begin{array}{c} b \end{array} \right) \end{array}$$

$S$  is ① random  
② not explicitly stored

# Frequent Elements (Heavy Hitters)

**Theorem 39.** The MISRA-GRIES algorithm is a deterministic one-pass algorithm which, for any given parameter  $\varepsilon \in (0, 1]$ , provides  $\hat{f}_1, \dots, \hat{f}_n$  of all element frequencies such that

$$f_j - \varepsilon m \leq \hat{f}_j \leq f_j, \quad j \in [n]$$

with space complexity  $s = O(\log(mn)/\varepsilon)$ . (In particular, it can be used to solve the MAJORITY problem in two passes.)

Remember Misra-Gries?

3 claims:

- ① MG is a sketching algorithm (yay!)
- ② ~~you will see that in tutorial~~
- ③ It is not a linear sketching algo

Not quite. See [AC lecture notes, p. 12, Exercise 1-3.]

Guarantee of MG: outputs  $\hat{f} \in \mathbb{R}^n$  (succinctly represented)  
s. t.  $\|\hat{f} - f\|_\infty \leq \varepsilon \|f\|_1$   
in the cash register model.

$$\|x\|_\infty = \max_{i \in [n]} |x_i|$$
$$\|f\|_1 = \sum_{i=1}^n |f_i|$$

## Frequent Elements (Heavy Hitters): $\ell_1, \ell_2$ , etc.

$$\|f\|_p = \left( \sum_{i=1}^n |f_i|^p \right)^{1/p}$$

$\|f\|_2 \approx$  measure  
of empirical  
variance  
of the stream

•  $\|f\|_0, \|f\|_1, \|f\|_2, \|f\|_\infty$  (fact)

$\uparrow$   
# of distinct  
elements  
( $F_0$ )

$$\|f\|_\infty \leq \|f\|_2 \leq \|f\|_1$$

" $\varepsilon \|f\|_2$  better guarantee than  $\varepsilon \|f\|_1$ "

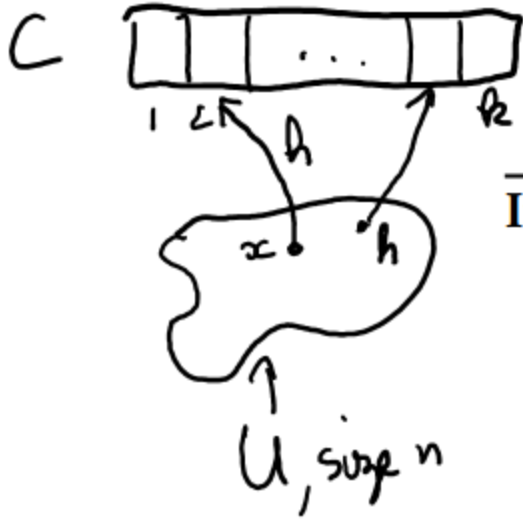
• In the rest: slight generalization:

$$S = \left( \underset{\substack{\uparrow \\ \text{demand}}}{(a_1, c_1)}, \dots, \underset{\substack{\uparrow \\ \text{\# of units}}}{(a_m, c_m)} \right)$$

$$|c_i| \leq B = O(1)$$

$$\rightarrow \|f\|_1 \leq m \cdot B$$

# Frequent Elements (Heavy Hitters): CountSketch (1/5)



**Input:** Parameters  $\epsilon, \delta \in (0, 1]$

- 1: Set  $k \leftarrow O(1/\epsilon^2)$ , and initialize an array  $C$  of size  $k$  to zero
- 2: Pick  $h: [n] \rightarrow [k]$  from a strongly universal hashing family
- 3: Pick  $g: [n] \rightarrow \{-1, 1\}$  from a strongly universal hashing family
- 4: **for all**  $1 \leq i \leq m$  **do**
- 5:     Get item  $a_i = (j, c) \in [n] \times \{-B, \dots, B\}$   $\triangleright$  Assume  $B = O(1)$
- 6:      $C[h(j)] \leftarrow C[h(j)] + c \cdot g(j)$

**Output:** On query  $j \in [n]$ , **return**  $\hat{f}_j \leftarrow g(j) \cdot C[h(j)]$

Turnstile model

$$\begin{array}{lll}
 a_1 = (25, 3) & h(25) = 2 & g(25) = -1 \\
 C[2] \leftarrow 0 + 3 \cdot (-1) = -3 \\
 a_{101} = (300, 1) & h(300) = 2 & g(300) = -1 \\
 C[2] \leftarrow -3 + 1 \cdot (-1) = -4
 \end{array}$$

$$\begin{array}{ll}
 a_{1005} = (25, -2) \\
 C[2] \leftarrow -4 + (-2) \cdot (-1) = -2
 \end{array}$$



## Frequent Elements (Heavy Hitters): CountSketch (2/5)

$$C[j] = \sum_{i: h(i)=j} \underbrace{g(i)}_{\pm 1} \cdot b_i$$

Example:  $\underbrace{g(25)}_{\hat{b}_{25}} C[2] = \underbrace{(g(25))^2}_{-1} b_{25} + \sum_{\substack{i \neq 25 \\ h(i)=2}} \underbrace{g(25)g(i)}_{\pm 1} b_i$

$$= b_{25} + \sum_{\substack{i \neq 25 \\ h(i)=2}} (\pm 1) b_i$$

noise with expectation 0

"Expect"

$$\hat{b}_i \approx b_i \pm O(\|b\|_2)$$


**Input:** Parameters  $\epsilon, \delta \in (0, 1]$

- 1: Set  $k \leftarrow O(1/\epsilon^2)$ , and initialize an array  $C$  of size  $k$  to zero
- 2: Pick  $h: [n] \rightarrow [k]$  from a strongly universal hashing family
- 3: Pick  $g: [n] \rightarrow \{-1, 1\}$  from a strongly universal hashing family
- 4: **for all**  $1 \leq i \leq m$  **do**
- 5:   Get item  $a_i = (j, c) \in [n] \times \{-B, \dots, B\}$   $\triangleright$  Assume  $B = O(1)$
- 6:    $C[h(j)] \leftarrow C[h(j)] + c \cdot g(j)$

**Output:** On query  $j \in [n]$ , **return**  $\hat{f}_j \leftarrow g(j) \cdot C[h(j)]$

Intuition

$$\sum_{i=1}^n \underbrace{x_i}_{\pm 1 \text{ uniformly}} \underbrace{a_i}_{\text{fixed vector}} \approx \mathcal{N}(0, \|a\|_2^2)$$

$b_{25} +$  



## Frequent Elements (Heavy Hitters): CountSketch (3/5)

Claim CountSketch is a sketching algorithm\*.

Proof. By looking at step "output".

Claim.  $\mathbb{E}[\hat{b}] = b$ .

Proof Fix  $j \in [n]$ .

$$\begin{aligned}\hat{b}_j &= g(j) C[h(j)] = g(j) \sum_{i: h(i)=h(j)} g(i) b_i \\ &= \underbrace{g(j) \cdot g(j)}_{=1} b_j + \sum_{\substack{i: i \neq j \\ h(i)=h(j)}} g(j) g(i) b_i\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\hat{b}_j] &= b_j + \sum_{i: i \neq j} \mathbb{E}[\mathbb{1}_{h(i)=h(j)} \underbrace{g(j)g(i)}_{=0}] b_i \\ &= b_j\end{aligned}$$

(pairwise indep\*)

**Input:** Parameters  $\epsilon, \delta \in (0, 1]$

- 1: Set  $k \leftarrow O(1/\epsilon^2)$ , and initialize an array  $C$  of size  $k$  to zero
- 2: Pick  $h: [n] \rightarrow [k]$  from a strongly universal hashing family
- 3: Pick  $g: [n] \rightarrow \{-1, 1\}$  from a strongly universal hashing family
- 4: **for all**  $1 \leq i \leq m$  **do**
- 5:   Get item  $a_i = (j, c) \in [n] \times \{-B, \dots, B\}$   $\triangleright$  Assume  $B = O(1)$
- 6:    $C[h(j)] \leftarrow C[h(j)] + c \cdot g(j)$

**Output:** On query  $j \in [n]$ , **return**  $\hat{f}_j \leftarrow g(j) \cdot C[h(j)]$

\* using the same hash functions  $g, h$  for all sketches.

## Frequent Elements (Heavy Hitters): CountSketch (3/5)

Claim  $\forall_j \quad \text{Var}[\hat{b}_j] \leq \frac{\|b\|_2^2}{k}$  ← size of  $C$

Proof  $\text{Var} \hat{b}_j = \mathbb{E}[\hat{b}_j^2] - \mathbb{E}[\hat{b}_j]^2 = \mathbb{E}[\hat{b}_j^2] - b_j^2$

$$\mathbb{E}[\hat{b}_j^2] = \mathbb{E}\left[\left(\sum_{j'} \mathbb{1}_{h(j')=h(j)} g(j')g(j) b_{j'}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j_1, j_2} \mathbb{1}_{h(j_1)=h(j_2)=h(j)} g(j_1)g(j_2) \cancel{g(j)}^2 b_{j_1}b_{j_2}\right]$$

$$= \sum_{j_1, j_2} \mathbb{E}[\mathbb{1}_{h(j_1)=h(j_2)=h(j)}] \underbrace{\mathbb{E}[g(j_1)g(j_2)]}_{= \mathbb{1}_{j_1=j_2}} b_{j_1}b_{j_2}$$

$$= \sum_{j_1} \underbrace{\mathbb{P}[h(j_1)=h(j)]}_{\substack{1 \text{ if } j_1=j \\ \frac{1}{k} \text{ o/w (pairwise indep.)}}} b_{j_1}^2$$

$$= b_j^2 + \frac{1}{k} \sum_{j_1 \neq j} b_{j_1}^2 \leq b_j^2 + \frac{1}{k} \sum_{j_1} b_{j_1}^2 = b_j^2 + \frac{\|b\|_2^2}{k}$$

From there,  
Chebyshev:

$$\begin{aligned} \mathbb{P}[|\hat{b}_j - b_j| > \varepsilon \|b\|_2] &\leq \frac{\|b\|_2^2/k}{\varepsilon^2 \|b\|_2^2} = \frac{1}{k\varepsilon^2} \\ &\leq \frac{1}{3} \end{aligned}$$

Input: Parameters  $\varepsilon, \delta \in (0, 1]$

- 1: Set  $k \leftarrow O(1/\varepsilon^2)$ , and initialize an array  $C$  of size  $k$  to zero
- 2: Pick  $h: [n] \rightarrow [k]$  from a strongly universal hashing family
- 3: Pick  $g: [n] \rightarrow \{-1, 1\}$  from a strongly universal hashing family
- 4: **for all**  $1 \leq i \leq m$  **do**
- 5:     Get item  $a_i = (j, c) \in [n] \times \{-B, \dots, B\}$  ▷ Assume  $B = O(1)$
- 6:      $C[h(j)] \leftarrow C[h(j)] + c \cdot g(j)$

Output: On query  $j \in [n]$ , **return**  $\hat{f}_j \leftarrow g(j) \cdot C[h(j)]$

## Frequent Elements (Heavy Hitters): CountSketch (5/5)

**Theorem 44.** *The (median trick version of the) COUNTSKETCH algorithm is a randomised one-pass sketching algorithm which, for any given parameters  $\epsilon, \delta \in (0, 1]$ , provides a (succinctly represented) estimate  $\hat{f}$  of frequency vector  $f$  of the stream such that, for every  $j \in [n]$*

$$\Pr \left[ \left| \hat{f}_j - f_j \right| \leq \epsilon \|f\|_2 \right] \geq 1 - \delta$$

*with space complexity*

$$s = O \left( \frac{\log(nm)}{\epsilon^2} \log \frac{1}{\delta} \right).$$

# Frequent Elements (Heavy Hitters): CountMinSketch (1/4)

⊛ Cash register model



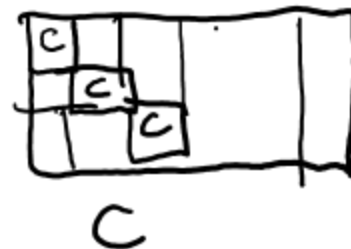
each row has  
its own  
 $h_t$

**Input:** Parameters  $\epsilon, \delta \in (0, 1]$

- 1: Set  $k \leftarrow O(1/\epsilon)$  and  $T \leftarrow O(\log(1/\delta))$ , and initialize a two-dimensional array  $C$  of size  $T \times k$  to zero
- 2: Pick  $h_1, \dots, h_T: [n] \rightarrow [k]$  independently from a strongly universal hashing family
- 3: **for all**  $1 \leq i \leq m$  **do**
- 4:     Get item  $a_i = (j, c) \in [n] \times \{0, \dots, B\}$  ▷ Assume  $B = O(1)$
- 5:     **for all**  $1 \leq t \leq T$  **do**     ⊛
- 6:          $C[t][h_t(j)] \leftarrow C[t][h_t(j)] + c$

**Output:** On query  $j \in [n]$ , **return**  $\hat{f}_j \leftarrow \min_{1 \leq t \leq T} C[t][h_t(j)]$

$$\begin{aligned} h_1(25) &= 1 \\ h_2(25) &= 2 \\ h_3(25) &= 3 \end{aligned}$$



"Bloom filter that counts"

## Frequent Elements (Heavy Hitters): CountMinSketch (2/4)

Claim.

Count Min Sketch is a sketching algo. \*

Claim.

(Cash register model)

Claim.

$$\hat{b}_j \geq b_j \quad \forall j.$$

$$\text{Space } O\left(\frac{\log(mn)}{\epsilon} \log\left(\frac{1}{\delta}\right)\right).$$

"Better" than CountSketch.

**Input:** Parameters  $\epsilon, \delta \in (0, 1]$

- 1: Set  $k \leftarrow O(1/\epsilon)$  and  $T \leftarrow O(\log(1/\delta))$ , and initialize a two-dimensional array  $C$  of size  $T \times k$  to zero
- 2: Pick  $h_1, \dots, h_T: [n] \rightarrow [k]$  independently from a strongly universal hashing family
- 3: **for all**  $1 \leq i \leq m$  **do**
- 4:   Get item  $a_i = (j, c) \in [n] \times \{0, \dots, B\}$  ▷ Assume  $B = O(1)$
- 5:   **for all**  $1 \leq t \leq T$  **do**
- 6:      $C[t][h_t(j)] \leftarrow C[t][h_t(j)] + c$

**Output:** On query  $j \in [n]$ , **return**  $\hat{f}_j \leftarrow \min_{1 \leq t \leq T} C[t][h_t(j)]$

\* using the same hash functions for each stream.

# Frequent Elements (Heavy Hitters): CountMinSketch (3/4)

**Input:** Parameters  $\epsilon, \delta \in (0, 1]$

- 1: Set  $k \leftarrow O(1/\epsilon)$  and  $T \leftarrow O(\log(1/\delta))$ , and initialize a two-dimensional array  $C$  of size  $T \times k$  to zero
- 2: Pick  $h_1, \dots, h_T: [n] \rightarrow [k]$  independently from a strongly universal hashing family
- 3: **for all**  $1 \leq i \leq m$  **do**
- 4:   Get item  $a_i = (j, c) \in [n] \times \{0, \dots, B\}$  ▷ Assume  $B = O(1)$
- 5:   **for all**  $1 \leq t \leq T$  **do**
- 6:      $C[t][h_t(j)] \leftarrow C[t][h_t(j)] + c$

**Output:** On query  $j \in [n]$ , **return**  $\hat{f}_j \leftarrow \min_{1 \leq t \leq T} C[t][h_t(j)]$

Claim

$$\forall j, \Pr[|\hat{f}_j - f_j| \geq \epsilon \|f\|_1] \leq \delta.$$

Proof:

$$\begin{aligned} \textcircled{1} \quad \hat{f}_j &= f_j + \text{"noise" (}\geq 0\text{) from collisions} \\ \hat{f}_j &= \min_{1 \leq t \leq T} C[t][h_t(j)] \\ &= \min_t \left( f_j + \sum_{\substack{j' \neq j \\ h_t(j') = h_t(j)}} f_{j'} \right) \end{aligned}$$

let's bound this

$$\textcircled{2} \quad \mathbb{E}[\hat{f}_j] = \frac{\|f\|_1}{k}$$

$$\begin{aligned} \textcircled{3} \text{ Marker. } \Pr[\hat{f}_j \geq \epsilon \|f\|_1] &\leq \frac{1}{\epsilon k} \\ &\leq \frac{1}{2} \end{aligned}$$

$$\begin{aligned} \textcircled{4} \text{ min bad } \Rightarrow \text{all } T \text{ are bad!} \\ \textcircled{*} \quad \Pr[\cdot] \leq \frac{1}{2^T} \end{aligned}$$

## Frequent Elements (Heavy Hitters): CountMinSketch (4/4)

**Theorem 45.** The COUNTMINSKETCH algorithm is a randomised one-pass sketching algorithm which, for any given parameters  $\epsilon, \delta \in (0, 1]$ , provides a (succinctly represented) estimate  $\hat{f}$  of frequency vector  $f$  of the stream such that, for every  $j \in [n]$

$$\Pr \left[ \left| \hat{f}_j - f_j \right| \leq \epsilon \|f\|_1 \right] \geq 1 - \delta$$

with space complexity

$$s = O \left( \frac{\log(nm)}{\epsilon} \log \frac{1}{\delta} \right).$$

(Moreover,  $\hat{f}_j$  is always an overestimate:  $\hat{f}_j \geq f_j$  for all  $j \in [n]$ .)



## Wait a minute...

This seems **strictly worse** than Misra-Gries!

- Randomised instead of deterministic!
- Uses more space!
- Also in the cash register model!
- Also an  $\ell_1$  guarantee!





## Wait a minute...

This seems **strictly worse** than Misra-Gries!

- Randomised instead of deterministic!
- Uses more space!
- Also in the cash register model!
- Also an  $\ell_1$  guarantee!

Yes, but:

- **Linear** sketch!
- Much **faster** per time step!
- Can be extended to the **strict turnstile** model!



## Recap

- CS, CMS, MG solve the same problem
- Different guarantees (CS/CMS, MG)
- Different types of sketches (CS, CMS / MG)
- Different space usage