

Count Min Algorithm : like a Bloom Filter

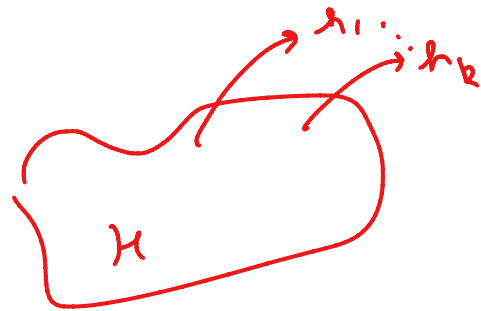
$\{1, \dots, n\}$



Stream  $S = \{a_1, \dots, a_m\}$ ,  $a_i \in \{1, \dots, n\} = [n]$

for  $\delta > 0$ ,  $\epsilon > 0$  (small values) Set  $t = \log_2\left(\frac{1}{\delta}\right)$ ,  $k = \frac{2}{\epsilon}$

choose  $t$  independent random hash functions  $h_1, \dots, h_t$ ;  $h_i: [n] \rightarrow [k]$



as maps these  
are deterministic

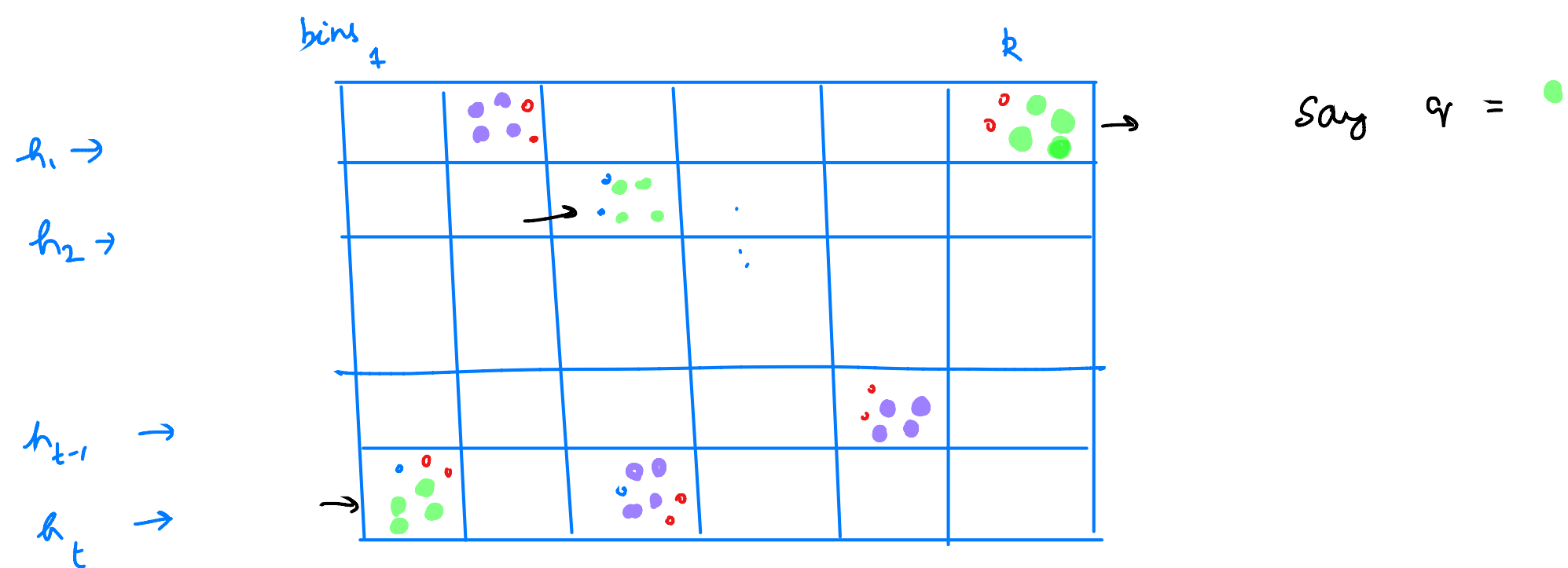
they map things far away  
from each other

Initialise  $t \times k$  counter matrix  $C$  s.t.  $C_{i,j} = 0$

for  $i \in [m]$  (sweep token wise)

do for  $j \in [t]$  (sweep hash function wise)

do  $C_{j, h_j(a_i)} = C_{j, h_j(a_i)} + 1$



For a query  $v \in [n]$ ,  $\hat{f}_v = \min_{j \in [t]} c_{j, h_j(v)}$

Count Algorithm : like a Bloom Filter

$\{1, \dots, n\}$



Stream  $S = \{a_1, \dots, a_m\}$ ,  $a_i \in \{1, \dots, n\} = [n]$

for  $\delta > 0$ ,  $\epsilon > 0$  (small values) Set  $t = \log_2\left(\frac{1}{\delta}\right)$ ,  $k = \frac{2}{\epsilon}$

choose  $t$  independent random hash functions  $h_1, \dots, h_t$ ;  $h_i : [n] \rightarrow [k]$

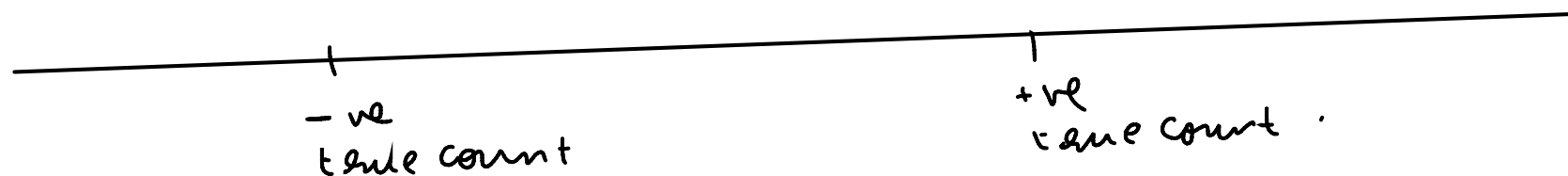
choose  $t$  "  $s_1, \dots, s_t$ ;  $s_i : [n] \rightarrow \{+1, -1\}$

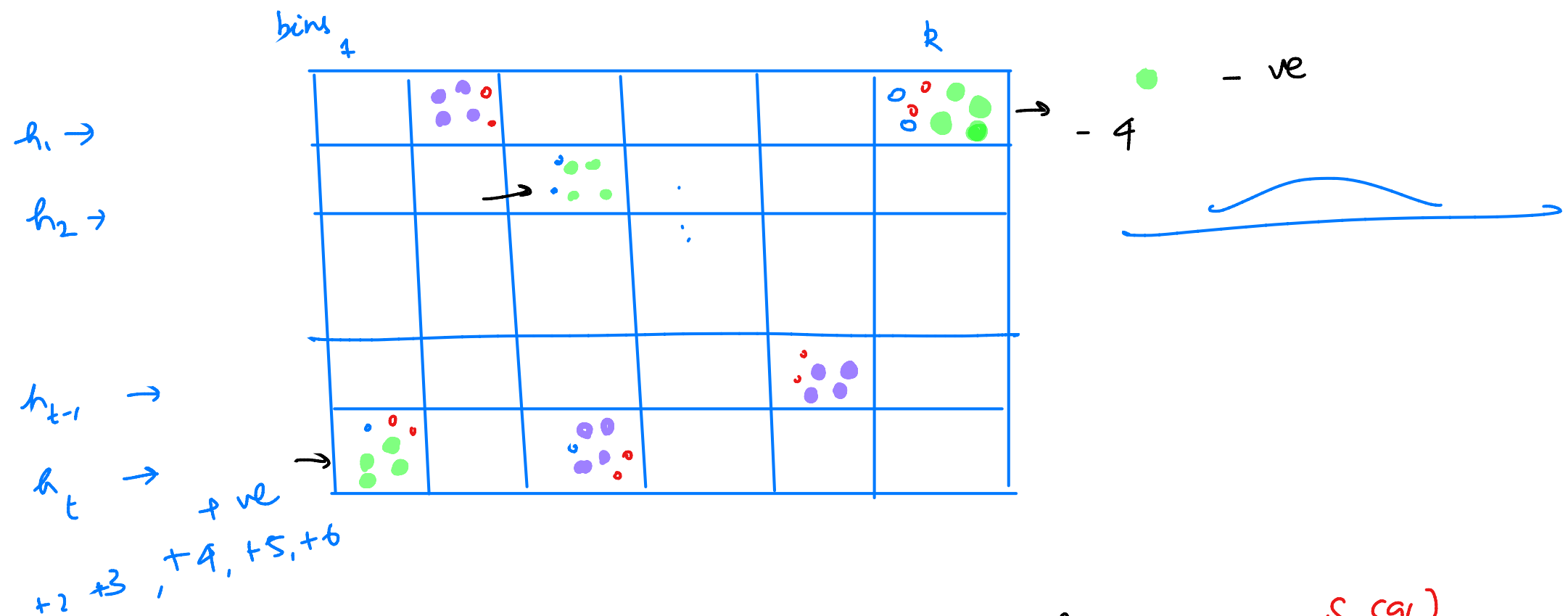
Initialise  $t \times k$  counter matrix  $C$  s.t.  $C_{i,j} = 0$

for  $i \in [m]$  (sweep token wise)

do for  $j \in [t]$  (sweep hash function wise)

do  $C_{j, h_j(a_i)} = C_{j, h_j(a_i)} + s_j(a_i)$





For a query  $q \in [n]$ ,  $\hat{f}_q = \underset{j}{\text{median}} \{c_{j, h_j(q)} \cdot s_j(q)\}$

Idea: frequent items = signal = same sign

Non-frequent items =  $+/-$  = cancels out each other

Actual = True  $\pm$  Spread

## Estimating Number of distinct elements

- $d = |\{i : f_i > 0\}|$

1, 2, 2, 2, 1, 3, 4, 1, 2,

- $\text{Zeros}(p) = \max \{i : 2^i \text{ divides } p\}$

$$\text{Zeros}(4) = 2, \quad \underbrace{100}$$

## Tide Mark Algorithm

- choose a random hash function

$$h: [n] \rightarrow [n]$$

- $z_j \leftarrow 0$

- Process (token  $j$ )

if  $\text{zeros}(h(j)) > z_j$  then  $z_j \leftarrow \text{zeros}(h(j))$

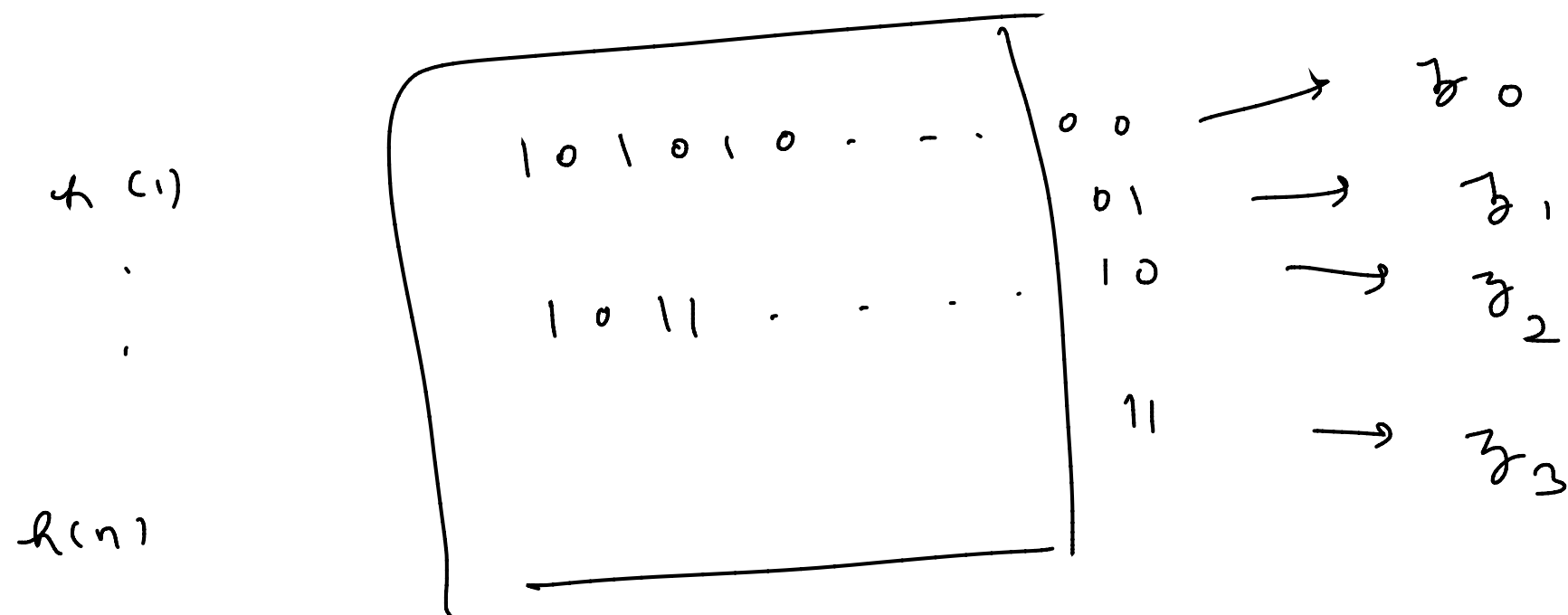
- output  $2^{z+\frac{1}{2}}$

$\neq \begin{matrix} z+1 \\ z+\frac{1}{2} \end{matrix}$

log log

\*  $z_1, \dots, z_m$  : counters

\* each item is directed to one of the 'm' counters.



\* estimate =  $m \cdot \sqrt{2} \cdot \frac{(z_1 + \dots + z_m)}{m}$

## Super log

\* same as log log

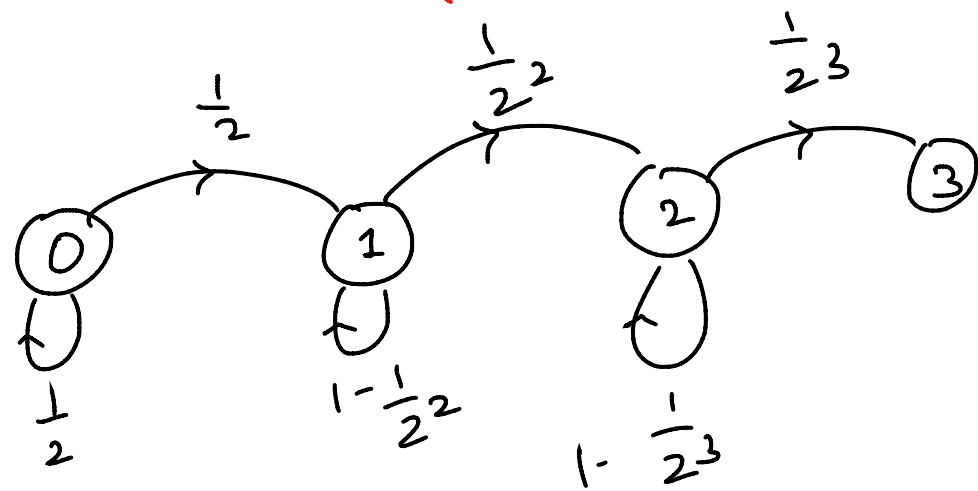
\* estimate =  $0.7 \cdot m \cdot \sqrt{2} \cdot 2^{\text{Avg (Smallest 70\% of } z_1, \dots, z_m \text{)}}$

## Hg per log log

\* estimate =  $m \cdot \sqrt{2} \cdot \text{harmonic mean } (2^{z_i})$



Idea: Simulate coin toss  $p = \text{prob of getting } 0$



- Probability of getting  $k$  consecutive heads with a fair coin is

$$\frac{1}{2^k}$$

- to see one such sample we need  $2^k$  trials.  
 $\text{count} = 2 \rightarrow \text{count} = 3$   
 3 consecutive 0s

$h(2) \rightarrow$  one coin toss sequence  $\rightarrow$

1	1	0	1	1	0	1	0	x
-	-	-	-	1	1	0	0	x
-	-	-	-	1	0	1		x
-	-	-	-	0	0	0	0	✓

•

$h(1) \rightarrow$  1 0 0 0 1 0 1  
 $h(2) \rightarrow$  1 0 1 1 1 0 1 0  
 .  
 .  
 .  
 .

$h(n) \rightarrow$

• to move from  $2^{y-1}$  to  $2^y$ , the counter has to see  $2^y$  elements (on an average)

• At  $2^y$  total distinct elements  
 $2 + 4 + \dots + 2^y = 2(1 + \dots + 2^{y-1})$   
 $= 2(2^y - 1) = 2^{y+1} \times \frac{1}{\sqrt{2}}$