

Misra-Gries Frequency Estimation

* Process σ

* Provide \hat{f}_j , $j \in \{1, \dots, n\}$

Init : $A \leftarrow \text{empty}$ (Balanced Binary Search tree)
 \uparrow approximate frequency estimator

Process (token j)

if $j \in \text{keys}(A)$

$$A[j] \leftarrow A[j] + 1$$

else if $|\text{keys}(A)| \leq k-1$ then

$$A[j] \leftarrow 1$$

else (if $|\text{keys}(A)| = k$)

for each $l \in \text{key}(A)$

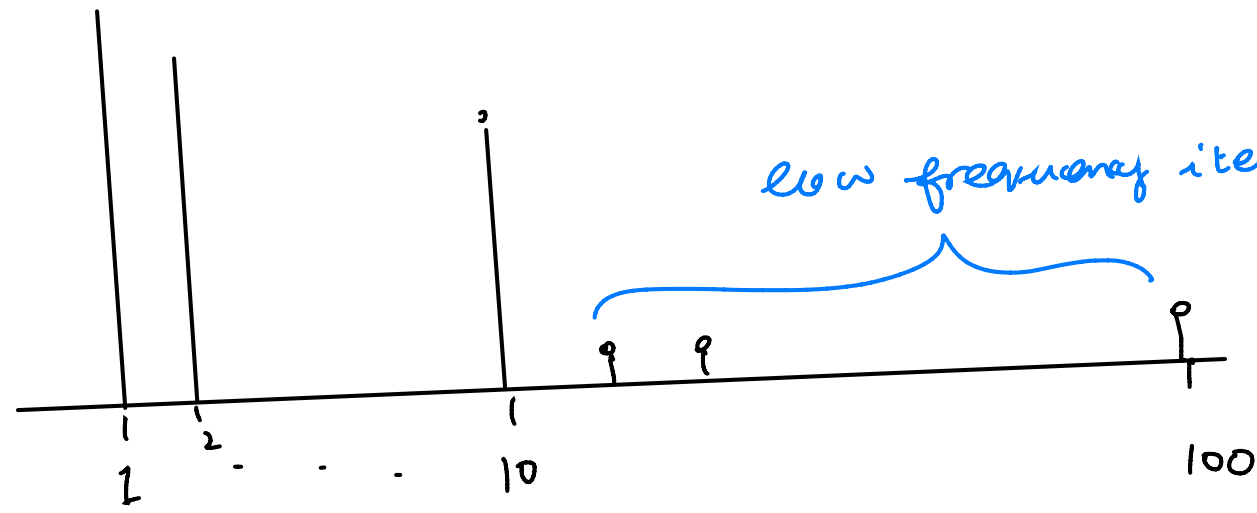
$$A[l] \leftarrow A[l] - 1$$

\rightarrow decrement step

if $A[l] = 0$, remove it from A

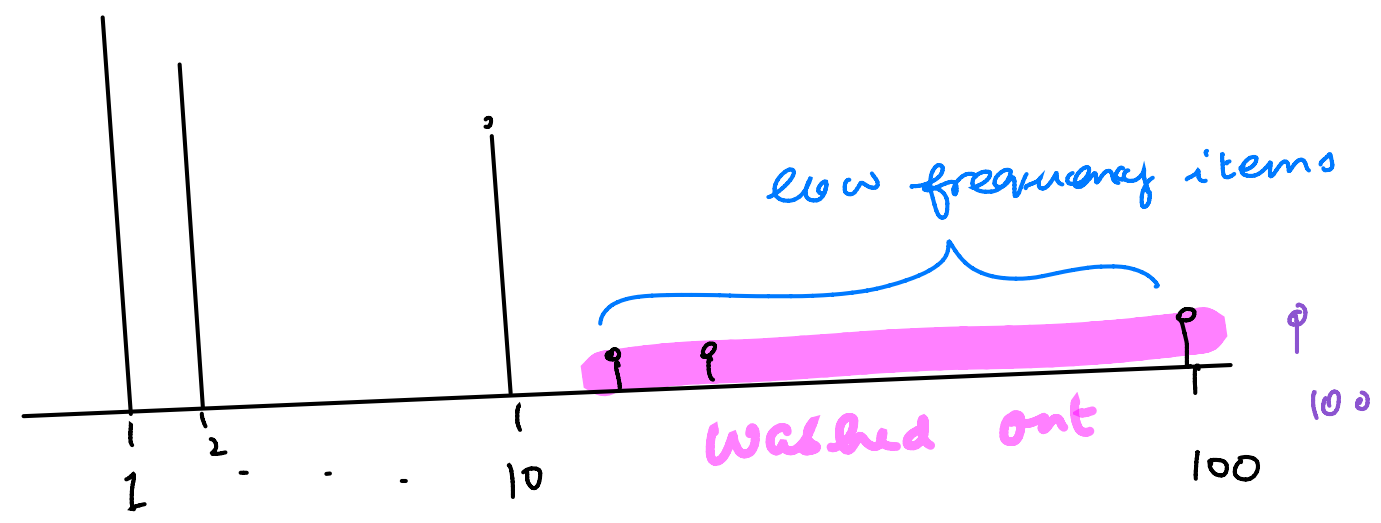


top 10, $n = 10^6$, $k = 100$



low frequency items have poor survival

top items are survival comparison



low frequency items have poor survival

top items are survival comparison

$n = 10^6$, top items are soap, Pen, mobile, laptop, pencil

a_1, \dots, a_m

Soap	Bag ₁	=	{ 2, 30, 31, 52	}	A[j]
Pen		=	{ 1, 20, 25	}	
mobile		=	{ 3, 18,	}	
laptop		=	{ 5,	}	
pencil		=	{ 4, . . .	}	
	Bag _k	=	{	}	

on the decrement step I think that the first items in the bag drop out

* Items come into bag once, if they leave they do not come back into the bag

* When decrement happens, k items leave at once, one from each bag.

\Rightarrow # decrements is at max $\frac{m}{k}$

$$f_j - \frac{m}{k} \leq \hat{f}_j \leq f_j$$

Typical scenario, $m = 10^6$, $k = 100$, top 10 are frequent
 \uparrow
search queries

$$f \sim \frac{m}{10} \sim 10^5, \quad \frac{m}{k} = \frac{10^6}{100} = 10^4$$

$$10^5 - 10^4 \leq \hat{f}_j \leq 10^5$$

$$f_f > \frac{m}{k} \Rightarrow \hat{f}_f > 0$$

But not the other way around \Rightarrow one can output infrequent items

Say, $k=2$, $m=5$, $S = \{a, b, c\}$

a, a, a, a, b

$$\hat{f}_1 = 4$$

$$\hat{f}_2 = 1$$

\Rightarrow item b is occurring more than $\frac{5}{2}$

a, a, b, c, a

$$\hat{f}_1 = 1, 2, 2, 1, 2$$

$$\hat{f}_2 = -, -, 1, -, -$$