

L11: Streaming : Frequent Items and Quantiles

Jeff M. Phillips

February 14, 2018

Streaming Model

Data $A = \langle a_1, a_2, \dots, a_i, \dots, a_n \rangle$

A_i treat as set $\{a_1, a_2, \dots, a_i\}$

Data Type

① $a_i \in [m] = \{1, 2, \dots, m\}$

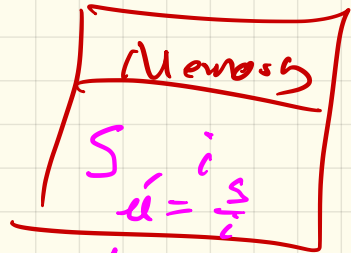
IP address

3-gram word
SNP

frequent items

$$S = \sum_{j=1}^i a_j$$

view a_i



much too small to store all data

② $a_i \in \mathbb{R}$ (scalar)
can be sorted

Apx CDF

Maintain statistics about A_i

Frequent Items

$$A = \langle a_1, a_2, \dots, a_n \rangle \quad a_i \in [m]$$

Mem Size $\ll n, m$

Approximate all $f_j = C \cdot \left(\underbrace{\log n}_{\text{counter}} + \underbrace{\log m}_{\text{label}} \right)$

frequency f_j $j \in [m]$

= # times $a_i = j$

$$= \left| \{ a_i \in A \mid a_i = j \} \right|$$

$$F_1 = \sum_{j=1}^m f_j \quad F_2 = \sqrt{\sum_{j=1}^m f_j^2}$$

$$F_0 = \sum_{j=1}^m f_j^0 = \# \text{ distinct items}$$

↳ 1 word

1 IP address

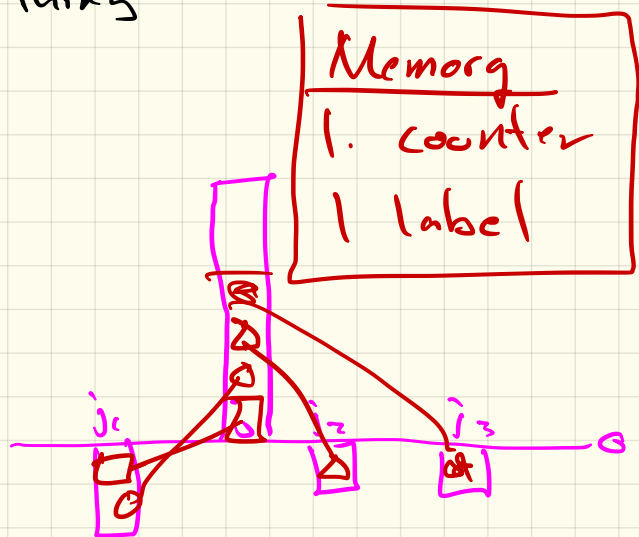
(counter, label)
 $O(\log m + \log n)$ space

↳ How many times saw IP address i .

MAJORITY

- if (some $f_j > \frac{n}{2}$), output j
else, output anything

```
Initialize  $L \neq C = 0$   
for  $i = 1$  to  $n$   
  Read  $a_i$   
  if ( $a_i = L$ )  $C = C + 1$   
  else  $C = C - 1$   
  if ( $C < 0$ ) :  $C = 1$  &  $L = a_i$   
end for  
return  $L$ 
```



Majority

Majority(A)

Set $c = 0$ and $\ell = \emptyset$

for $i = 1$ **to** m **do**

if $(a_i = \ell)$ **then**

$c = c + 1$

else

$c = c - 1$

if $(c < 0)$ **then**

$c = 1, \ell = a_i$

return ℓ

Misra - Gries

For any $j \in [m]$ Return $S(j) = \hat{f}_j$ $\epsilon = \text{error} \in [0, 1]$

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

$$k = \frac{1}{\epsilon} \Rightarrow \frac{n}{k} = n\epsilon$$

Initialize:

for $i=1$ to n

- If ($a_i = \text{some } L_j$) $C_j = C_j + 1$
- else (if some $C_j = 0$) $C_j = 1$ $L_j = a_i$
- else

Decrement all counters

for $l=1$ to $k-1$ $C_l = C_l - 1$

can occur at most n/k times

Memory
($k-1$) counters + labels

end for

Return $S [C = C_1, C_2, \dots, C_{k-1} \mid L = L_1, L_2, \dots, L_{k-1}]$

Misra-Gries

counter array $C : C[1], C[2], \dots, C[k - 1]$

location array $L : L[1], L[2], \dots, L[k - 1]$

Misra-Gries(A)

Set all $C[i] = 0$ and all $L[i] = \emptyset$

for $i = 1$ **to** m **do**

if ($a_i = L[j]$) **then**

$C[j] = C[j] + 1$

else

if (some $C[j] = 0$) **then**

 Set $L[j] = a_i$ & $C[j] = 1$

else

for $j \in [k - 1]$ **do** $C[j] = C[j] - 1$

return C, L

Quantiles

$$A = \langle a_1, a_2, \dots, a_n \rangle$$

$$q_i \in \mathbb{R}$$

↑ requires $\log m$ space

$$\text{rank}_A(v) = |\{a_i \in A \mid a_i \leq v\}|$$

$v \in \mathbb{R}$

$$\text{cdf}(v) = \frac{\text{rank}_A(v)}{n}$$

Summary Q_A s.t.

$$\forall v \in \mathbb{R}$$

$$\left| Q_A(v) - \frac{\text{rank}_A(v)}{n} \right| < \epsilon$$

error $\in [0, 1]$

Memory

Space $(\log m) \cdot \frac{1}{\epsilon} \log \log \frac{1}{\epsilon}$

$$k = \frac{1}{\epsilon}$$

Set $\frac{1}{k}$ values sorted order $Q [g_1 \leq g_2 \leq \dots \leq g_k]$

Merge $Q, Q' \Rightarrow g_1 \leq g'_1 \leq g_2 \leq g'_2 \leq g_3 \leq g'_3 \leq g_4 \leq g'_4 \leq \dots$

Frugal Median

Let $v \in \mathbb{R}$ desired so $\frac{\text{rank}_A(v)}{n} = \frac{1}{2}$
 \uparrow
median $\rightarrow \hat{g}$

s.t. \hat{g} close to v .

$$\frac{\text{rank}_A(\hat{g})}{n} \in \left[\frac{1}{2} - \epsilon, \frac{1}{2} + \epsilon \right]$$

Memory
1 label

Set $l = a_1$

for $i = 2$ to n

if $(a_i > l)$

$l = l + 1$

else $(a_i < l)$

$l = l - 1$

return l

improve estimation
increment

Frugal Median

Frugal Median(A)

Set $l = 0$.

for $i = 1$ **to** m **do**

if $(a_i > l)$ **then**

$l \leftarrow l + 1$.

if $(a_i < l)$ **then**

$l \leftarrow l - 1$.

return l .

What if estimator
is for
25%-quantile?

Frugal Quantile

Frugal Quantile(A, ϕ)

e.g. $\phi = 0.75$

Set $l = 0$.

for $i = 1$ **to** m **do**

$r = \text{Unif}(0, 1)$ (at random)

if ($a_i > l$ **and** $r > 1 - \phi$) **then**

$l \leftarrow l + 1$.

if ($a_i < l$ **and** $r > \phi$) **then**

$l \leftarrow l - 1$.

return l .