

# L14: Streaming : Frequent Items and Quantiles

Jeff M. Phillips

February 25, 2025

Data is Big! Too big to fit on computer  
or in memory

1. Parallelism More computers  
MapReduce

2. Sampling very large  $X$  of data  
 $S \sim X$   $S = \{s_1, s_2, \dots, s_n\}$   
 $n \ll |X|$

3. Streaming  $X = \langle x_1, x_2, x_3, \dots, x_i, \dots, x_n \rangle$   
Read each  $x_i$  once, but not store  
Maintain small space summary.

Data

$$A = \langle a_1, a_2, a_3, \dots, a_i, \dots, a_n \rangle$$

ex.  $a_i \in \mathbb{R}$

$$A_i = \langle a_1, a_2, \dots, a_i \rangle$$

mean ( $A_i$ )

$$\frac{1}{i} \sum_{j=1}^i a_j = v_i$$

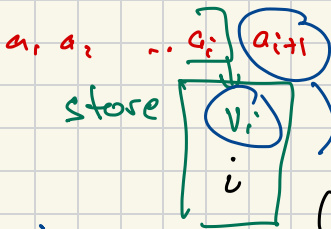
sum  $S_i = \sum_{j=1}^i a_j$

$$v_i = \frac{S_i}{i}$$

$i = \# \text{ steps}$

$\oplus a_{i+1}$

$$\begin{aligned} S_{i+1} &= S_i + a_{i+1} \\ i+1 &= i + 1 \end{aligned}$$



$$\frac{(v_i \cdot i) + a_{i+1}}{i+1}$$

$$\frac{1}{i} \sum_{j=1}^i (a_j - E[A_i])^2 = \text{Var}[A_i]$$

$$v_{i+1} = \frac{S_{i+1}}{i+1}$$

Variance  $[A_i]$

$$\begin{aligned} \frac{1}{i} \sum_{j=1}^i (a_j - E[A_i])^2 &= \frac{1}{i} \sum_{j=1}^i a_j^2 - \left( \frac{1}{i} \sum_{j=1}^i a_j \right)^2 \\ &= E[A_i^2] - (E[A_i])^2 \end{aligned}$$

so we

$i$
$S_i$
$Q_i$

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

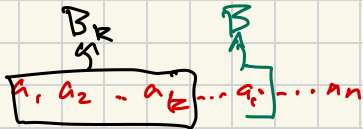
$$Q_i = \sum_{j=1}^i (a_j)^2$$

$$\text{Var}(A_i) = \frac{Q_i}{i} - \left( \frac{S_i}{i} \right)^2$$

# Reservoir Sampling $(a_1, a_2, \dots, a_n)$ Stream

Goal Maintain Random Sample  $B \sim A$   
without replacement  $|B| = k$

1. Keep first  $k$  items  $B = \{a_1, a_2, \dots, a_k\}$
2. for  $j = k+1 \dots n$   
keep  $a_j$  in  $B$  w/p  $\frac{k}{j}$   
↳ boot a random  $b \in B$  from  $B$ .  
o.w. keep  $B$  the same



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

$a_i \in [m]$

$[m] = \{ \text{Addresses} \}$

$[m] = \{ \text{Keys} \}$

both  $m$  and  $n$  too big.

label  $j \in [m]$   $\log m$  bits

counter  $i \in [n]$   $\log n$  bits

approx

freq

$f_j$

$|f_j - f_j^*| \leq \epsilon n$

frequency  $f_j = |\{ a_i \in A \mid a_i = j \}|$   
 $j \in [m]$

Heavy-Hitter

$j$  s.t.  $f_j > \phi n$

$\phi = 0.1$

$\epsilon = 0.01$

ok  $f_j > \phi n - \epsilon n$

# MAJORITY

$$A = (a_1, \dots, a_n)$$
$$a_i \in \{m\}$$

if (some  $f_j > \frac{n}{2}$ )  $\rightarrow$  output  $j$

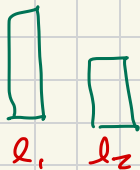
else output anything  $j \in \{m\}$

[ counter  $\in \{n\}$  | label  $j \in \{m\}$  ]  
e l

if ( $a_i = l$ )  $c = c + 1$

else  $c = c - 1$

| if ( $c < 0$ )  
 $l = a_i$ ;  $c = 1$



# Majority

Majority(A)

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

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

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

$c = c + 1$

**else**

$c = c - 1$

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

$c = 1, \ell = a_i$

**return**  $\ell$



# Frequency Approximation

# Misra-Gries

$\forall_j \in [m]$   $f_j$   $k-1$  counters, s.t.  $f_j - \frac{n}{k} \leq f_j \leq f_j$   $k-1$  labels  $k = \frac{1}{\epsilon}$

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

Labels  $l_1, l_2, \dots, l_{k-1} = L$   
counters  $c_1, c_2, \dots, c_{k-1}$

for  $(a_i \in A)$

if  $(a_i \in L)$   $a_i = l_j$

$c_j = c_j + 1$

else  $(a_i \notin L)$

if (some  $c_j \leq 0$ )

$l_j = a_i$

$c_j = 1$

no labels  
unused.

decrease

all counter

by 1  $\forall_j c_j = c_j - 1$



# 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**  $n$  **do**

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

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

**else**

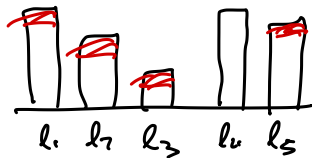
**if** (some  $C[j] = 0$ ) **then**

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

**else**

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

**return**  $C, L$



how many  
total  
decrements?

$$\hookrightarrow \frac{n}{k}$$



# 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**  $l$ .

# Frugal Quantile

Frugal Quantile( $A, \phi$ )

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**  $l$ .

## Frequent Itemsets : Apriori

$$T_1 = \{1, 2, 3, 4, 5\}$$

$$T_2 = \{2, 6, 7, 9\}$$

$$T_3 = \{1, 3, 5, 6\}$$

$$T_4 = \{2, 6, 9\}$$

$$T_5 = \{7, 8\}$$

$$T_6 = \{1, 2, 6\}$$

$$T_7 = \{0, 3, 5, 6\}$$

$$T_8 = \{0, 2, 4\}$$

$$T_9 = \{2, 4\}$$

$$T_{10} = \{6, 7, 9\}$$

$$T_{11} = \{3, 6, 9\}$$

$$T_{12} = \{6, 7, 8\}$$