

L11: Streaming : Frequent Items and Quantiles

Jeff M. Phillips

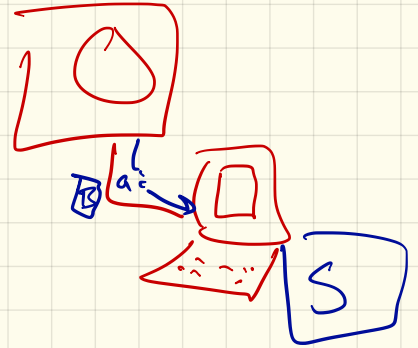
February 20, 2019

Big Data

• Sampling

• Streaming : One Pass

• MapReduce
Distributed



Streaming

m, n too big

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

$a_i \in [n]$ Domain

OR size

counter = $O(\log m)$ bits

label = $O(\log n)$ bits

$[n] = \mathbb{IP}$ address

= words in dictionaries

k-gram

hash tables (hash function)
 $O(\log n + \log m)$

Streaming Model

$$A = \langle a_1, a_2, \dots, a_j, \dots, a_m \rangle$$

$$a_j \in [n] \leftarrow \text{Domain}$$

m, n Very Large

frequencies

Space: $(\log m + \log n)$
counter label

$$f_j = |\{a_i \in A \mid a_i = j\}|$$

$$F_0 = \sum_i f_i^0 = \# \text{ distinct elements}$$

$$F_1 = \sum_j f_j^1 = m = \# \text{ elements}$$

Hyperloglog

$$F_2 = \sqrt{\sum_j f_j^2} = \text{join size}$$

MAJORITY

Is one IP address on more than half of all packets?

Is some $f_j > m/2$?

If so, which one?

Then report j s.t. $f_j > m/2$.

If not, guess.

If not, return any $j \in [n]$.

Majority

Majority(A)

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

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

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

$c = c + 1$

increment

else

$c = c - 1$

decrement

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

$c = 1, l = a_i$

return l

Heavy Hitters / Frequent Items

Report all $f_i \geq m/k$

For all $i \in [n]$ have $f_i \leftarrow \text{approx}$

$$f_i - \frac{m}{k} \leq f_i \leq f_i + \frac{m}{k}$$

MG

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

$$\hookrightarrow \frac{m}{k} = \epsilon m$$

$$\epsilon = \frac{1}{10} \text{ error}$$

$$\epsilon = 0.01 \Rightarrow 1\% \text{ error}$$

$\hookrightarrow k = 100$

Misra - Gries Algo

Cache: $k-1$ counters

$k-1$ labels $\leftarrow k$ guesses

$C[1]$ $C[2]$... $C[k]$

$L[1]$ $L[2]$... $L[k]$

• if $a_j = L[j]$ $C[j]++$

• else and some $C[j] = 0$ \leftarrow empty counter
then $L[j] = a_j$ $C[j] = 1$

• else Decrement all counters
 $\forall j' \in [k] \quad C[j']--$

Total #
Decrements
 $\leq \frac{m}{k}$

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

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

total # Decrements

$$m = 1600$$

$$k = 100$$

$$\hat{f}_j = \cancel{150} 140$$

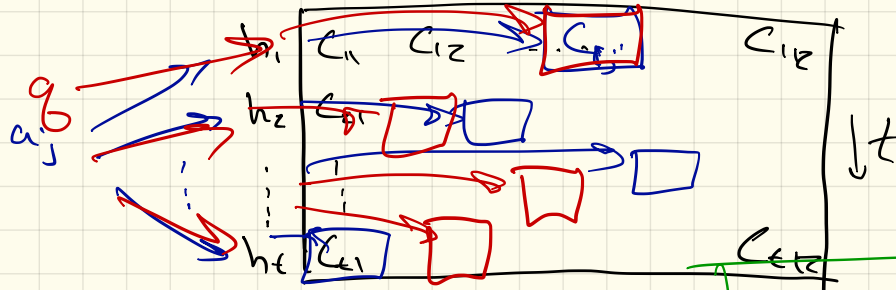
$$Q. \text{ is } f_j > 0.15 \cdot m?$$

Count - Min Sketch

turnstyle model

$k \cdot t$ counters

t hash functions
 $h: [n] \rightarrow [k]$



$\tilde{O}(\log n + \log m)$
space

$k = \frac{n}{\epsilon}$ $t = \log \frac{1}{\delta}$
 prob. failure

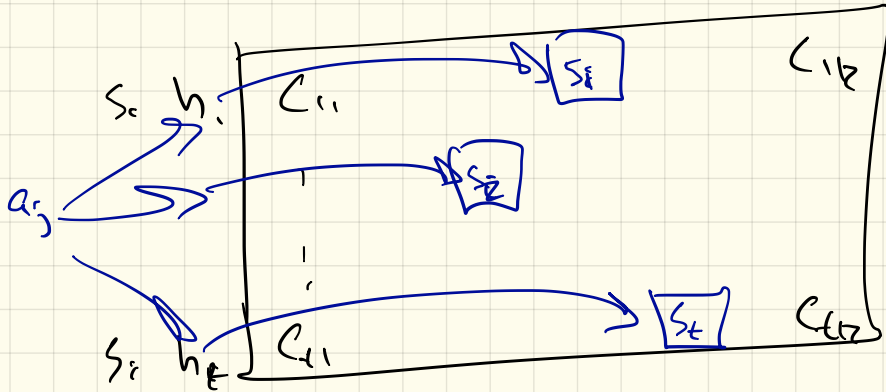
for $a_j \in A$
 for $i = 1 \dots t$
 $C_{i, h_i(j)} ++$

linear sketch

Query $g \in [n]$
 $\hat{f}_g = \min_i C_{i, h_i(g)}$

$f_g \leq \hat{f}_g \leq f_g + \epsilon m$ w.p. $> 1 - \delta$

Count Sketch



t counters

t hash functions

t random sign hashes

$$S_i: [n] \rightarrow \{-1, +1\}$$

for $a_j \in A$

for i in $1 \dots t$

$$C_{i, h_i}(a_j) = C_{i, h_i}(a_j) + S_i(j)$$

Query g
 $\hat{f}_g = \text{median}_i (S_i(g) \cdot C_{i, h_i}(g))$

$$t = O\left(\frac{1}{\epsilon^2}\right)$$

$$\leftarrow \log^{2/5}$$

$$E[\hat{f}_g] = f_g$$

$$|\hat{f}_g - f_g| \leq \epsilon \cdot F_2$$

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, 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 .

Frequent Itemsets : Apriori

$$A = \{T_1, T_2, \dots, T_m\}$$

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

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

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

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

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

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

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

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

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

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

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

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

