

L11: Streaming : Frequent Items and Quantiles

Jeff M. Phillips

February 12, 2020

Big Data

size $|X| = n$
very large

too big to fit on 1 computer.

- Parallelism: More computers
MapReduce

- Sampling: $P \sim X \stackrel{\text{iid}}{\sim} \mathcal{U}$ $|P| \ll |X|$
data set $P \sim \mathcal{U}$

- Streaming: $X = \langle x_1, x_2, \dots, x_i, \dots, x_n \rangle$
Read data in one pass
Maintain small space summary

Data $A = \langle a_1, a_2, \dots, a_n \rangle$
 $a_i \in \mathbb{R}$

mean $(A_i) = \frac{S_i}{i}$ maintain

variance $(A_i) = \frac{Q_i}{i} + \left(\frac{S_i}{i}\right)^2$

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

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

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

Reservoir Sampling

Maintain Random Sample

$B \subset A$

w/o replacement
uniform

$$|B| = k$$

1. Keep first k $B = A_k$

2. for $j = k+1$ to n

with prob $\frac{k}{j}$

→ replace some $b_i^r \in B$
w/ a_j

otherwise

keep B

ϵ - error

$$\hookrightarrow k = \frac{1}{\epsilon^2}$$

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

$$= 1, 2, \dots, m$$

n too large

m too large

IP addresses

k-grams

label $j \in [m]$ $\frac{\log m}{\log n}$ bits

counter value $\in \{1, \dots, n\}$ $\frac{\log n}{\log n}$ bits

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

$j \in [m]$

$$\text{approx freq } \hat{f}_j : \left| f_j - \hat{f}_j \right| \leq \epsilon n$$

$$\text{Is any } j \in [m] \text{ have } f_j > \epsilon n$$

$$\text{or } \hat{f}_j > \epsilon n - \epsilon n$$

MAJORITY

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

if some $f_j > \frac{n}{2} \rightarrow \text{output } j$
else output anything.

1 counter, 1 label

if $(a_i = \text{label})$
counter $c = c + 1$

else $c = c - 1$

if $c < 0$
 $l = a_i$

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 ℓ

Misra - Greis

Freq. Apx

$k-1$ counters, $k-1$ labels

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

decrements \uparrow (pointing to $f_j - \frac{n}{k}$)

apx freq \uparrow (pointing to \hat{f}_j)

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

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

for all $j \in [m]$

if j not in set $L = \{l_1, l_2, \dots, l_{k-1}\}$

$\hat{f}_j = 0$

large counts \uparrow (pointing to L)

storing $\hat{f}_j = 1$ counter
1 label j

counters $C = \{c_1, \dots, c_{k-1}\}$

for $(a_i \in A)$

if $a_i \in L$ \leftarrow numbers on label l_j

$$c_j = c_j + 1$$

else $(a_i \notin L)$

for $j \in [1 \dots k-1]$ $c_j = c_j - 1$

if $(c_j \in C \text{ has } c_j \leq 0)$
 $\hat{f}_j = a_i, c_j = 1$

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** m **do**

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

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

else

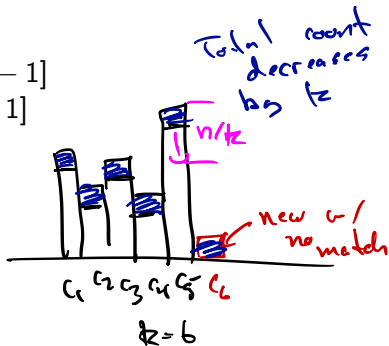
if (some $C[j] = 0$) **then**

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

else

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

return C, L



How many times?

n/k times

Streaming Median

$$A = \langle a_1, \dots, a_n \rangle \quad a_i \in \mathbb{R}$$

Maintain median (A)

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

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

quantile estimate

$$\text{for } v \in \mathbb{R} \quad \left| \Phi_A(v) - \frac{\text{rank}_A(v)}{n} \right| \leq \epsilon$$

$$\frac{1}{\epsilon} \log \log \frac{1}{\epsilon}$$

space

$$\frac{1}{\epsilon} \log \frac{1}{\epsilon}$$

Frugal Median

Assume A iid
order random

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 .