# L15 : CountMin Sketch

## (and friends)

Jeff M. Phillips | Data Mining

# Stream

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

$$a_i \in [m] \quad \leftarrow \text{domain}$$

Compute Statistics
on $A$

$n$ too large

$\log n \quad \leftarrow \text{counter}$

- small space

$m$ too large

$\log m \quad \leftarrow \text{label}$

- one pass

$$S \; \text{stetch}(A)$$

frequency $\quad j \in [m]$

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

$$F_1 = \sum_j f_j = n$$

$$F_2 = \sqrt{\sum_j f_j^2} \qquad F_2 \ll F_1$$

$$F_0 = \text{number of distinct elements in } A.$$

# Refresh Frequency Approximation

$$\forall_j \in [m] \quad \rightarrow \quad \hat{f}_j \qquad \text{so} \quad |f_j - \hat{f}_j| \leq \varepsilon n = \varepsilon F_1$$

$$\text{MG:} \quad f_j - \varepsilon n \leq \hat{f}_j \leq f_j \qquad\qquad \leq \varepsilon F_2 \quad \text{(Count}$$
$$\text{Sketch)}$$

$$\text{(CountMin} \qquad f_j \leq \hat{f}_j \leq f_j + \varepsilon n$$

$$\text{w.p } 1-\delta$$
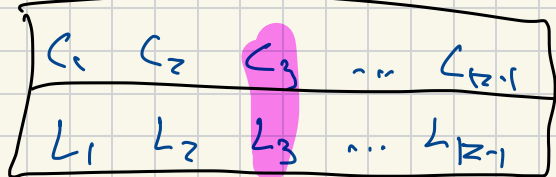
$$k = 1/\varepsilon$$

$$k-1 \quad \text{counters}$$
$$k-1 \quad \text{labels}$$

## Sketch Data Structure

Data Structure $S(A)$
- update $S(A)$ w/ $a_i$
- query $S(A)(j) \Rightarrow \hat{f}_j$

| $C_1$ | $C_2$ | $C_3$ | ... | $C_{k-1}$ |
|---|---|---|---|---|
| $L_1$ | $L_2$ | $L_3$ | ... | $L_{k-1}$ |

$$a_i \in A \qquad \text{elseif } (C_j = 0) \qquad \xrightarrow{\text{else}} \text{Decrem}$$
$$\text{if } a_i = L_j \quad | \quad L_j = a_i \quad C_j = 1 \quad | \quad \text{all counters}$$

# Count Min Sketch

$t \cdot (k = 2/\epsilon)$ counters

$t = \log_2 \frac{1}{\delta}$ hash functions          $h_j : [m] \to [k]$

$a$



$k$ columns

$h_1(a)$  $C_{11}$   $C_{12}$         $C_{1,h_1(a)}$        $C_{1k}$

$h_2(a)$              $C_{2,h_2(a)}$

$t$ rows          $C_{jr}$                    $C_{ji}$

$h_t$   $C_{t1}$   $C_{t2}$                                $C_{tk}$

randomness

$h_j \sim \mathcal{H}$

uniform

$\forall_j \in [t]$
$C_{j, h_j(q)} \geq f_q$

$q \in [m]$
$\Rightarrow \hat{f_q}$

## insert $(a_i \in A)$

for $j = 1$ to $t$

$C_{j, h_j(a_i)} ++$

## query $(q)$ $\Rightarrow \hat{f_q}$

$$\hat{f_q} = \min_{j \in [t]} C_{j, h_j(q)}$$

$f_{\hat{g}} \leq \hat{f}_{\hat{g}}$ each counter $C_{j, h_j(\hat{g})}$ includes count of $f_{\hat{g}}$

$\hat{f}_{\hat{g}} \leq f_{\hat{g}} + \boxed{w_{\hat{g}}} \qquad w \leq \varepsilon n = \xi F_1 = \xi \cdot \sum_j^1 f_j \qquad \boxed{t = \frac{2/\varepsilon}{}}$

---

Some $s \in [m]$ $\quad w_{\hat{g}}(s) = Y_s = \begin{cases} f_s & \text{if } h_j(s) = h_j(\hat{g}) \quad w.p. \; 1/k \\ 0 & o.w. \end{cases}$

Some $\hat{j} \in [t]$

total overcount $\quad w = \sum_{s \in [m]} w(s) = \sum_{s \in [m]} Y_s = X$

$E[X] = E\left[\sum_s Y_s\right] = \sum_s E[Y_s] = \sum_s f_s / k = \frac{1}{k} \sum_s f_s = \frac{1}{k} F_1 = \frac{\varepsilon n}{2}$

Markov Ineq.

R.V. $X > 0$ $\qquad \alpha = E[X] \cdot 2 \qquad \Longrightarrow \quad Pr[X > \varepsilon n] \leq \frac{E[X]}{E[X] \cdot 2} = \frac{1}{2}$

$Pr[X > \alpha] < \frac{E[X]}{\alpha} \qquad = \varepsilon n$

## $t$ hash functions $\qquad t = \log_2(1/\delta)$

1 hash $\quad h_j \sim \mathcal{H}$

$$\Pr\left[\omega_j = X > \varepsilon n\right] \leq \frac{1}{2}$$

$$\Pr\left[\text{all hash } h_j \quad \text{has} \quad \omega_j > \varepsilon n\right] = \left(\frac{1}{2}\right)^t$$

$t$ independent hash functions

$$\left(\frac{1}{2}\right)^t = \left(\frac{1}{2}\right)^{\log_2(1/\delta)} = 2^{-\log_2(1/\delta)} = 2^{\log_2(\delta)}$$

$$\delta = \frac{1}{2^{10}} = \frac{1}{1024} \qquad\qquad = \delta = \frac{1}{32}$$

$$t = 10 \qquad\qquad\qquad \Rightarrow t = 5$$

Compare MG vs Count Min

|  | MG | Count Min |
|---|---|---|
| Space | $1/\varepsilon$ counters + labels | $\frac{2}{\varepsilon} \cdot \log_2 \frac{1}{\delta}$ counters + $\log \frac{1}{\delta}$ hashfxn |
|  | Deterministic | Randomized w.p. $1-\delta$ |
| Bias: | under count | over count |
|  | heavy hitters | most query guesses or HAG |
| Deletions | X | can handle deletions (linear sketch) |

# Count Sketch

Unbiased $\varepsilon F_2$

$k \quad \text{columns}$

$t$ hash functions

$h_1$

$h_2(q)$

$h_j$

$h_k$

$C_{i,j}$

$h_j \sim \mathcal{H} \qquad h_j : [m] \to [k]$

Sign $\qquad S_j \sim \mathcal{S} \qquad S_j : [m] \to \{-1, +1\}$

query $(q) \qquad q \in [m]$

$\hat{f}_q = \text{median}\left( C_{j, h_j(q)} \right)$

$k = \dfrac{4}{\varepsilon^2}$

$t = 2 \log_2 \dfrac{1}{\delta}$

insert $(a_i)$

for $j = 1$ to $t$

$C_{j, h_j(a_i)} \mathrel{+}= S_j(a_i)$

$E[C_{j,i}] = 0$

$\left| \hat{f}_{\cdots} - f_{\cdots} \right| \leq \varepsilon F_2$

$E[\hat{f}_i] = f_i$